# Memory in the Age of AI Agents: A Survey
## Forms, Functions and Dynamics

**Yuyang Hu**[†], **Shichun Liu**[†], **Yanwei Yue**[†], **Guibin Zhang**[†]⬡, **Boyang Liu**, **Fangyi Zhu**, **Jiahang Lin**, **Honglin Guo**, **Shihan Dou**, **Zhiheng Xi**, **Senjie Jin**, **Jiejun Tan**, **Yanbin Yin**, **Jiongnan Liu**, **Zeyu Zhang**, **Zhongxiang Sun**, **Yutao Zhu**, **Hao Sun**, **Boci Peng**, **Zhenrong Cheng**, **Xuanbo Fan**, **Jiaxin Guo**, **Xinlei Yu**, **Zhenhong Zhou**, **Zewen Hu**, **Jiahao Huo**, **Junhao Wang**, **Yuwei Niu**, **Yu Wang**, **Zhenfei Yin**, **Xiaobin Hu**, **Yue Liao**, **Qiankun Li**, **Kun Wang**, **Wangchunshu Zhou**, **Yixin Liu**, **Dawei Cheng**, **Qi Zhang**, **Tao Gui**[‡], **Shirui Pan**, **Yan Zhang**[‡], **Philip Torr**, **Zhicheng Dou**[‡], **Ji-Rong Wen**, **Xuanjing Huang**[‡], **Yu-Gang Jiang**, **Shuicheng Yan**[‡]

[†]Core Contributors with Names Listed Alphabetically. ⬡ Project Organizer. [‡]Core Supervisors.

**Affiliations**: National University of Singapore, Renmin University of China, Fudan University, Peking University, Nanyang Technological University, Tongji University, University of California San Diego, Hong Kong University of Science and Technology (Guangzhou), Griffith University, Georgia Institute of Technology, OPPO, Oxford University

Memory has emerged, and will continue to remain, a core capability of foundation model-based agents. It underpins long-horizon reasoning, continual adaptation, and effective interaction with complex environments. As research on agent memory rapidly expands and attracts unprecedented attention, the field has also become increasingly fragmented. Existing works that fall under the umbrella of agent memory often differ substantially in their motivations, implementations, assumptions, and evaluation protocols, while the proliferation of loosely defined memory terminologies has further obscured conceptual clarity. Traditional taxonomies such as long/short-term memory have proven insufficient to capture the diversity and dynamics of contemporary agent memory systems. This survey aims to provide an up-to-date and comprehensive landscape of current agent memory research. We begin by clearly delineating the scope of agent memory and distinguishing it from related concepts such as LLM memory, retrieval augmented generation (RAG), and context engineering. We then examine agent memory through the unified lenses of **forms**, **functions**, and **dynamics**. From the perspective of forms, we identify three dominant realizations of agent memory, namely *token-level*, *parametric*, and *latent memory*. From the perspective of functions, we move beyond coarse temporal categorizations and propose a finer-grained taxonomy that distinguishes *factual*, *experiential*, and *working memory*. From the perspective of dynamics, we analyze how memory is formed, evolved, and retrieved over time as agents interact with their environments. To support empirical research and practical development, we compile a comprehensive summary of representative benchmarks and open source memory frameworks. Beyond consolidation, we articulate a forward-looking perspective on emerging research frontiers, including automation-oriented memory design, the deep integration of reinforcement learning with memory systems, multimodal memory, shared memory for multi-agent systems, and trustworthiness issues. We hope this survey serves not only as a reference for existing work, but also as a conceptual foundation for rethinking memory as a first-class primitive in the design of future agentic intelligence.
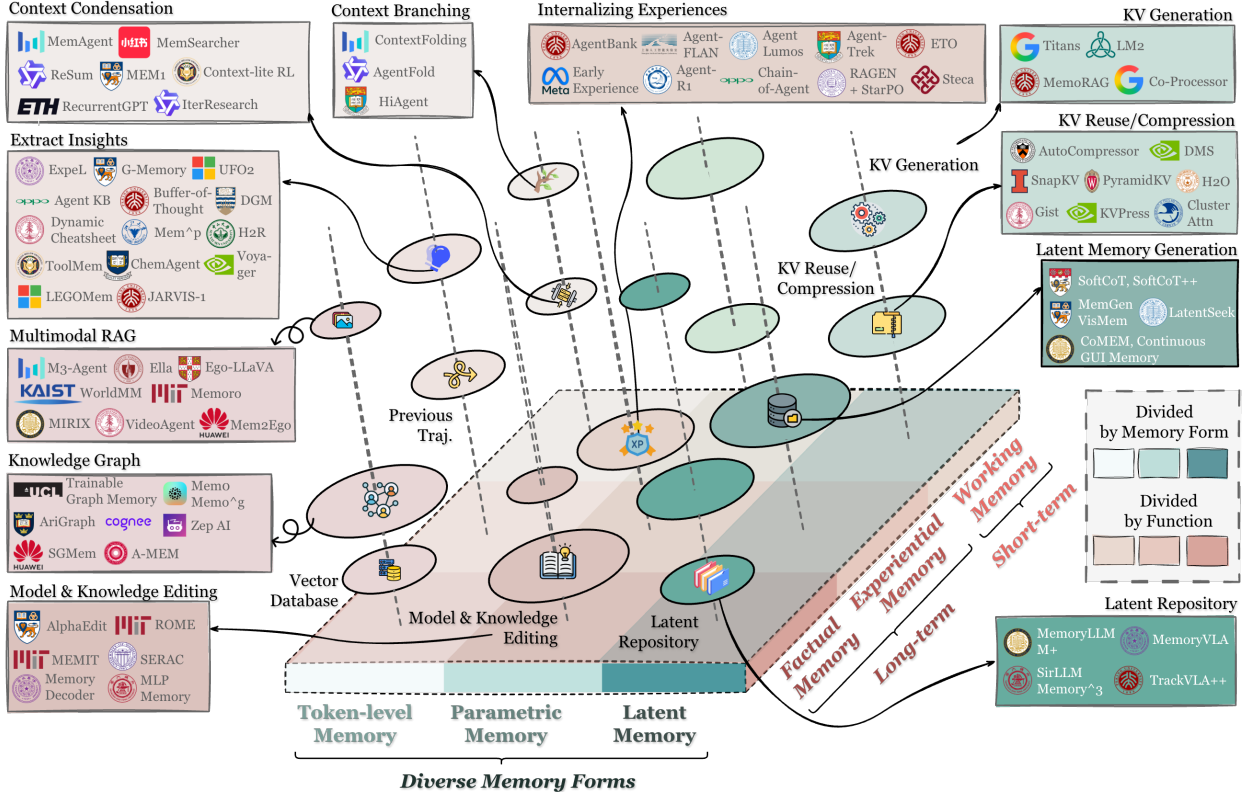
✉ **Main Contact:** guibinz@u.nus.edu, yuyang.hu@ruc.edu.cn, liusc24@m.fudan.edu.cn, ywyue25@stu.pku.edu.cn

⯁ **Github:** https://github.com/Shichun-Liu/Agent-Memory-Paper-List

## 1 Introduction

The past two years have witnessed the overwhelming evolution of increasingly capable large language models

---

Note: If you identify your own or other papers relevant to this survey that have not been discussed (we apologize for any such omissions due to the rapidly expanding literature), please feel free to contact us via email or raise an issue on GitHub.

**Figure 1** Overview of agent memory organized by the unified taxonomy of *forms* (Section 3), *functions* (Section 4), and *dynamics* (Section 5). The diagram positions memory artifacts by their dominant form and primary function. It further maps representative systems into this taxonomy to provide a consolidated landscape.

(LLMs) into powerful AI agents (Matarazzo and Torlone, 2025; Minaee et al., 2025; Luo et al., 2025). These foundation-model-powered agents have demonstrated remarkable progress across diverse domains such as deep research (Xu and Peng, 2025; Zhang et al., 2025o), software engineering (Wang et al., 2024i), and scientific discovery (Wei et al., 2025c), continuously advancing the trajectory toward artifacial general interlligence (AGI) (Fang et al., 2025a; Durante et al., 2024). Although early conceptions of "agents" were highly heterogeneous, a growing consensus has since emerged within the community: beyond a pure LLM backbone, an agent is typically equipped with capabilities such as *reasoning*, *planning*, *perception*, *memory*, and *tool-use*. Some of these abilities, such as reasoning and tool-use, have been largely internalized within model parameters through reinforcement learning (Wang et al., 2025l; Qu et al., 2025b), while some still depend heavily on external agentic scaffolds. Together, these components transform LLMs from static conditional generators into learnable policies that can interact with diverse external environments and adaptively evolve over time (Zhang et al., 2025f).

Among these agentic faculties, *memory* stands out as a cornerstone, explicitly enabling the transformation of static LLMs, whose parameters cannot be rapidly updated, into adaptive agents capable of continual adaptation through environmental interaction (Zhang et al., 2025r; Wu et al., 2025g). From an application perspective, numerous domains demand agents with proactive memory management rather than ephemeral, forgetful behaviors: personalized chatbots (Chhikara et al., 2025; Li et al., 2025b), recommender systems (Liu et al., 2025b), social simulations (Park et al., 2023; Yang et al., 2025), and financial investigations (Zhang et al., 2024) all rely on the agent's ability to process, store, and manage historical information. From a developmental standpoint, one of the defining aspirations of AGI research is to endow agents with the capacity for continual evolution through environment interactions (Hendrycks et al., 2025), a capability fundamentally grounded in agent memory.

**Agent Memory Needs A New Taxonomy**   Given the growing significance and community attention surrounding agent memory systems, it has become both timely and necessary to provide an updated perspective on contemporary agent memory research. The motivation for a new taxonomy and survey is twofold: ❶ **Limitations of Existing Taxonomies:** While several recent surveys have provided valuable and comprehensive overviews of agent memory (Zhang et al., 2025r; Wu et al., 2025g), their taxonomies were developed prior to a number of rapid methodological advances and therefore do not fully reflect the current breadth and complexity of the research landscape. For example, emerging directions in 2025, such as memory frameworks that distill reusable tools from past experiences (Qiu et al., 2025a,c; Zhao et al., 2025c), or memory-augmented test-time scaling methods (Zhang et al., 2025g; Suzgun et al., 2025), remain underrepresented in earlier classification schemes. ❷ **Conceptual Fragmentation:** With the explosive growth of memory-related studies, the concept itself has become increasingly expansive and fragmented. Researchers often find that papers claiming to study "agent memory" differ drastically in implementation, objectives, and underlying assumptions. The proliferation of diverse terminologies (declarative, episodic, semantic, parametric memory, etc.) further obscures conceptual clarity, highlighting the urgent need for a coherent taxonomy that can unify these emerging concepts.

Therefore, this paper seeks to establish a systematic framework that reconciles existing definitions, bridges emerging trends, and elucidates the foundational principles of memory in agentic systems. Specifically, this survey aims to address the following key questions:

> **Key Questions**
>
> ❶ How is *agent memory* defined, and how does it relate to related concepts such as LLM memory, retrieval-augmented generation (RAG), and context engineering?
>
> ❷ **Forms:** What architectural or representational forms can agent memory take?
>
> ❸ **Functions:** Why is agent memory needed, and what roles or purposes does it serve?
>
> ❹ **Dynamics:** How does agent memory operate, adapt, and evolve over time?
>
> ❺ What are the promising frontiers for advancing agent memory research?

To address question ❶, we first provide formal definitions for LLM-based agents and agent memory systems in Section 2, and present a detailed comparison between agent memory and related concepts such as LLM memory, RAG, and context engineering. Following the "Forms–Functions–Dynamics" triangle, we offer a structured overview of agent memory. Question ❷ examines the architectural forms of memory, which we discuss in Section 3, highlighting three mainstream implementations: token-level, parametric, and latent memory. Question ❸ concerns the functional roles of memory, addressed in Section 4, where we distinguish between *factual memory*, which records knowledge from agents' interactions with users and the environment; *experiential memory*, which incrementally enhances the agent's problem-solving capabilities through task execution; and *working memory*, which manages workspace information during individual task instances. Question ❹ focuses on the lifecycle and operational dynamics of agent memory, which we present sequentially in terms of memory formulation, retrieval, and evolution.

After surveying existing research through the lenses of "Forms–Functions–Dynamics," we further provide our perspectives and insights on agent memory research. To facilitate knowledge sharing and future development, we first summarize key benchmarks and framework resources in Section 6. Building upon this foundation, we then address question ❺ by exploring several emerging yet underdeveloped research frontiers in Section 7, including automation-oriented memory design, the integration of reinforcement learning (RL), multimodal memory, shared memory for multi-agent systems, and trustworthy issues.

**Contributions**   The contributions of this survey can be summarized as follows: (1) We present an up-to-date and multidimensional taxonomy of agent memory from the perspective of "forms–functions–dynamics," offering a structured lens through which to understand current developments in the field. (2) We provide an in-depth discussion on the suitability and interplay of different memory forms and functional purposes, offering insights into how various memory types can be effectively aligned with distinct agentic objectives. (3) We investigate emerging and promising research directions in agent memory, thereby outlining future

opportunities and guiding pathways for advancement. (4) We compile a comprehensive collection of resources, including benchmarks and open-source frameworks, to support both researchers and practitioners in further exploration of agent memory systems.

**Outline of the Survey.** The remainder of this survey is organized as follows. Section 2 formalizes LLM-based agents and agent memory systems, and clarifies their relationships with related concepts. Section 3, Section 4, and Section 5 respectively examine the forms, functions, and dynamics of agent memory. Section 6 summarizes representative benchmarks and framework resources. Section 7 discusses emerging research frontiers and future directions. Finally, we conclude the survey with a summary of key insights in Section 8.

## 2 Preliminaries: Formalizing Agents and Memory

LLM agents increasingly serve as the decision-making core of interactive systems that operate over time, manipulate external tools, and coordinate with humans or other agents. To study memory in such settings, we begin by formalizing LLM-based agent systems in a manner that encompasses both single-agent and multi-agent configurations. We then formalize the memory system coupled to the agent's decision process through read/write interactions, enabling a unified treatment of memory phenomena that arise both *within* a task (inside-trial / short-term memory) and *across* tasks (cross-trial / long-term memory).

### 2.1 LLM-based Agent Systems

**Agents and Environment** Let $\mathcal{I} = \{1, \ldots, N\}$ denote the index set of agents, where $N = 1$ corresponds to the single-agent case (e.g., ReAct), and $N > 1$ represents multi-agent settings such as debate (Li et al., 2024c) or planner–executor architectures (Wan et al., 2025). The environment is characterized by a state space $\mathcal{S}$. At each time step $t$, the environment evolves according to a controlled stochastic transition model

$$s_{t+1} \sim \Psi(s_{t+1} \mid s_t, a_t),$$

where $a_t$ denotes the action executed at time $t$. In multi-agent systems, this abstraction allows for either sequential decision-making (where a single agent acts at each step) or implicit coordination through environment-mediated effects. Each agent $i \in \mathcal{I}$ receives an observation

$$o_t^i = O_i(s_t, h_t^i, \mathcal{Q}),$$

where $h_t^i$ denotes the portion of the interaction history visible to agent $i$. This history may include previous messages, intermediate tool outputs, partial reasoning traces, shared workspace states, or other agents' contributions, depending on the system design. $\mathcal{Q}$ denotes the task specification, such as a user instruction, goal description, or external constraints, which is treated as fixed within a task unless otherwise specified.

**Action Space** A distinguishing feature of LLM-based agents is the heterogeneity of their action space. Rather than restricting actions to plain text generation, agents may operate over a multimodal and semantically structured action space, including:

- **Natural-language generation**, such as producing intermediate reasoning, explanations, responses, or instructions (Li et al., 2023b; Wu et al., 2024b; Hong et al., 2024; Qian et al., 2024).

- **Tool invocation actions**, which call external APIs, search engines, calculators, databases, simulators, or code execution environments (Qin et al., 2025; Li et al., 2025g; Zhou et al., 2023c, 2024c).

- **Planning actions**, which explicitly output task decompositions, execution plans, or subgoal specifications to guide later behavior (CAMEL-AI, 2025; Liu et al., 2025f; Pan et al., 2024).

- **Environment-control actions**, where the agent directly manipulates the external environment (e.g., navigation in embodied settings (Shridhar et al., 2021; Wang et al., 2022a), editing a software repository (Jimenez et al., 2024; Aleithan et al., 2024), or modifying a shared memory buffer).

- **Communication actions**, enabling collaboration or negotiation with other agents through structured messages (Marro et al., 2024).

These actions, though diverse in semantics, are unified by the fact that they are produced through an autoregressive LLM backbone conditioned on a contextual input. Formally, each agent $i$ follows a policy

$$a_t = \pi_i(o_t^i, m_t^i, \mathcal{Q}),$$

where $m_t^i$ is a memory-derived signal defined in Section 2.2. The policy may internally generate multi-step reasoning chains, latent deliberation, or scratchpad computations prior to emitting an executable action; such internal processes are abstracted away and not explicitly modeled.

**Interaction Process and Trajectories**   A full execution of the system induces a trajectory

$$\tau = (s_0, o_0, a_0, s_1, o_1, a_1, \ldots, s_T),$$

where $T$ is determined by task termination conditions or system-specific stopping criteria. At each step, the trajectory reflects the interleaving of (i) environment observation, (ii) optional memory retrieval, (iii) LLM-based computation, and (iv) action execution that drives the next state transition.

This formulation captures a broad class of agentic systems, ranging from a single agent solving reasoning tasks with tool augmentation to teams of role-specialized agents collaboratively developing software (Qian et al., 2024; Wang et al., 2025k) or conducting scientific inquiry (Weng et al., 2025). We next formalize the memory systems that integrate into this agent loop.

## 2.2  Agent Memory Systems

While an LLM-based agent interacts with an environment, its instantaneous observation $o_t^i$ is often insufficient for effective decision-making. Agents therefore rely on additional information derived from prior interactions, both within the current task and across previously completed tasks. We formalize this capability through a unified *agent memory system*, represented as an evolving memory state

$$\mathcal{M}_t \in \mathbb{M},$$

where $\mathbb{M}$ denotes the space of admissible memory configurations. No specific internal structure is imposed on $\mathcal{M}_t$; it may take the form of a text buffer, key–value store, vector database, graph structure, or any hybrid representation. At the beginning of a task, $\mathcal{M}_t$ may already contain information distilled from prior trajectories (cross-trial memory). During task execution, new information accumulates and functions as short-term, task-specific memory. Both roles are supported within a single memory container, with temporal distinctions emerging from usage patterns rather than architectural separation.

**Memory Lifecycle: Formation, Evolution, and Retrieval.**   The dynamics of the memory system are characterized by three conceptual operators.

**Memory Formation.**  At time step $t$, the agent produces informational artifacts $\phi_t$, which may include tool outputs, reasoning traces, partial plans, self-evaluations, or environmental feedback. A formation operator

$$\mathcal{M}_{t+1}^{\text{form}} = F(\mathcal{M}_t, \phi_t)$$

selectively transforms these artifacts into memory candidates, extracting information with potential future utility rather than storing the entire interaction history verbatim.

**Memory Evolution.**  Formed memory candidates are integrated into the existing memory base through an evolution operator

$$\mathcal{M}_{t+1} = E(\mathcal{M}_{t+1}^{\text{form}}),$$

which may consolidate redundant entries (Zhao et al., 2024), resolve conflicts (Rasmussen et al., 2025; Li et al., 2025k), discard low-utility information (Wang et al., 2025q), or restructure memory for efficient retrieval. The resulting memory state persists across subsequent decision steps and tasks.

**Memory Retrieval.** When selecting an action, agent $i$ retrieves a context-dependent memory signal

$$m_t^i = R(\mathcal{M}_t, o_t^i, \mathcal{Q}),$$

where $R$ denotes a retrieval operator that constructs a task-aware query and returns relevant memory content. The retrieved signal $m_t^i$ is formatted for direct consumption by the LLM policy, for example as a sequence of textual snippets or a structured summary.

**Temporal Roles Within the Agent Loop.** Although memory is represented as a unified state $\mathcal{M}_t$, the three lifecycle operators (formation $F$, evolution $E$, and retrieval $R$) need not be invoked at every time step. Instead, different memory effects arise from distinct temporal invocation patterns. For instance, some systems perform retrieval only once at task initialization,

$$m_t^i = \begin{cases} R(\mathcal{M}_0, o_0^i, \mathcal{Q}), & t = 0, \\ \bot, & t > 0, \end{cases}$$

where $\bot$ denotes null retrieval strategy. Others may retrieve memory intermittently or continuously based on contextual triggers. Similarly, memory formation may range from minimal accumulation of raw observations,

$$\mathcal{M}_{t+1}^{\text{form}} = \mathcal{M}_t \cup \{o_t^i\},$$

to sophisticated extraction and refinement of reusable patterns or abstractions. Thus, *inside a task*, short-term memory effects may arise from lightweight logging just as in Yao et al. (2023b); Chen et al. (2023a) or from more elaborate iterative refinement (Hu et al., 2025a); *across tasks*, long-term memory may be updated episodically at task boundaries or continuously throughout operation. Short-term and long-term memory phenomena therefore emerge not from discrete architectural modules but from the temporal patterns with which formation, evolution, and retrieval are engaged.

**Memory–Agent Coupling.** The interaction between memory and the agent's decision process is similarly flexible. In general, the agent policy is written as

$$a_t = \pi_i(o_t^i, m_t^i, \mathcal{Q}),$$

where the retrieved memory signal $m_t^i$ may be present or absent depending on the retrieval schedule. When retrieval is disabled at a given step, $m_t^i$ can be treated as a distinguished null input.

Consequently, the overall agent loop consists of observing the environment, optionally retrieving memory, computing an action, receiving feedback, and optionally updating memory through formation and evolution. Different agent implementations instantiate different subsets of these operations at different temporal frequencies, giving rise to memory systems that range from passive buffers to actively evolving knowledge bases.
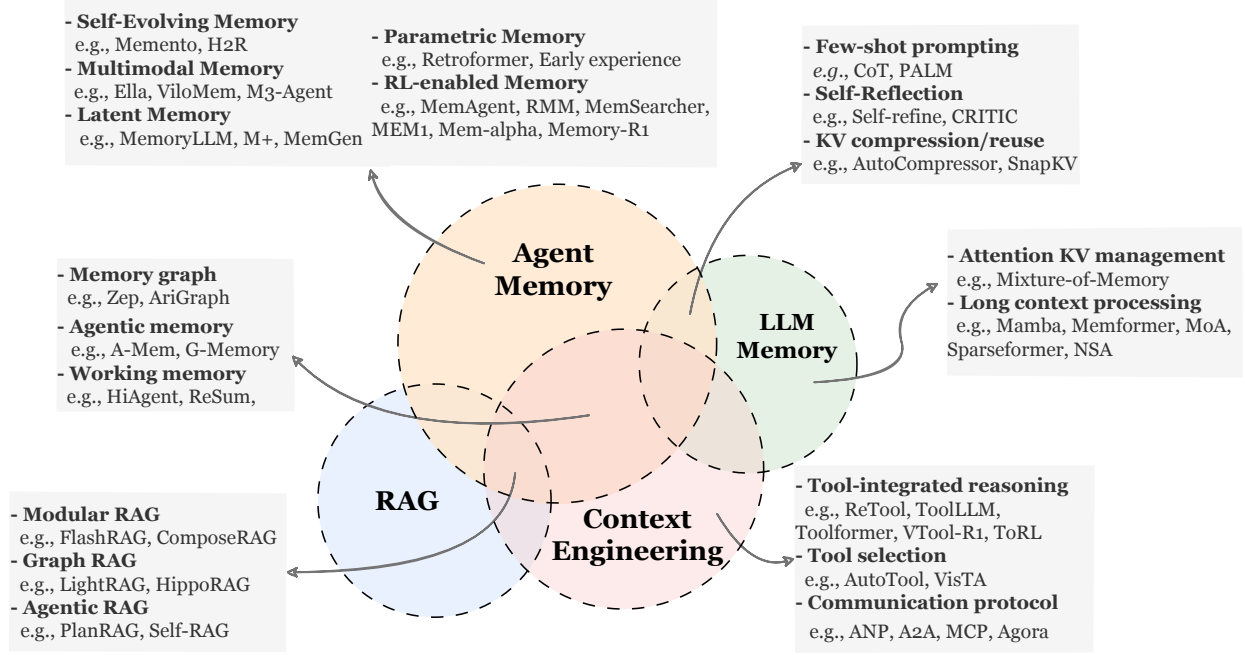
## 2.3 Comparing Agent Memory with Other Key Concepts

Despite the growing interest in agentic systems endowed with memory, the community's understanding of what constitutes *agent memory* remains fragmented. In practice, researchers and practitioners often conflate agent memory with related constructs such as LLM memory (Wu et al., 2025g), retrieval-augmented generation (RAG) (Gao et al., 2024), and context engineering (Mei et al., 2025). Although these concepts are intrinsically connected by their involvement in how information is managed and utilized in LLM-driven systems, they differ in scope, temporal characteristics, and functional roles.

These overlapping yet distinct notions have led to ambiguity in the literature and practice. To clarify these distinctions and situate agent memory within this broader landscape, we examine how agent memory *relates to*, and *diverges from*, LLM memory, RAG, and context engineering in the subsequent subsubsections. Figure 2 visually illustrates the commonalities and distinctions among these fields through a Venn diagram.

### 2.3.1 Agent Memory vs. LLM Memory

At a high level, *agent memory* almost fully subsumes what has traditionally been referred to as *LLM memory*. Since 2023, many works describing themselves as "LLM memory mechanisms" (Zhong et al., 2024; Packer et al., 2023a; Wang et al., 2023b) are more appropriately interpreted, under contemporary terminology, as

- **Self-Evolving Memory**
  e.g., Memento, H2R
- **Multimodal Memory**
  e.g., Ella, ViloMem, M3-Agent
- **Latent Memory**
  e.g., MemoryLLM, M+, MemGen

- **Parametric Memory**
  e.g., Retroformer, Early experience
- **RL-enabled Memory**
  e.g., MemAgent, RMM, MemSearcher, MEM1, Mem-alpha, Memory-R1

- **Few-shot prompting**
  *e.g.*, CoT, PALM
- **Self-Reflection**
  e.g., Self-refine, CRITIC
- **KV compression/reuse**
  e.g., AutoCompressor, SnapKV

**Agent Memory**

**LLM Memory**

- **Attention KV management**
  e.g., Mixture-of-Memory
- **Long context processing**
  e.g., Mamba, Memformer, MoA, Sparseformer, NSA

- **Memory graph**
  e.g., Zep, AriGraph
- **Agentic memory**
  e.g., A-Mem, G-Memory
- **Working memory**
  e.g., HiAgent, ReSum,

**RAG**

**Context Engineering**

- **Modular RAG**
  e.g., FlashRAG, ComposeRAG
- **Graph RAG**
  e.g., LightRAG, HippoRAG
- **Agentic RAG**
  e.g., PlanRAG, Self-RAG

- **Tool-integrated reasoning**
  e.g., ReTool, ToolLLM, Toolformer, VTool-R1, ToRL
- **Tool selection**
  e.g., AutoTool, VisTA
- **Communication protocol**
  e.g., ANP, A2A, MCP, Agora

**Figure 2** Conceptual comparison of **Agent Memory** with **LLM Memory**, **RAG**, and **Context Engineering**. The diagram illustrates shared technical implementations (e.g., KV reuse, graph retrieval) while highlighting fundamental distinctions: unlike the architectural optimizations of LLM Memory, the static knowledge access of RAG, or the transient resource management of Context Engineering, Agent Memory is uniquely characterized by its focus on maintaining a persistent and self-evolving cognitive state that integrates factual knowledge and experience. The listed categories and examples are illustrative rather than strictly parallel, serving as representative reference points to clarify conceptual relationships rather than to define a rigid taxonomy.

early instances of agent memory. This reinterpretation arises from the historical ambiguity surrounding the very notion of an "LLM agent." During 2023–2024, the community had no stable or coherent definition: in some cases, prompting an LLM to call a calculator already sufficed to qualify the system as an agent (Wu et al., 2024c); in other cases, agency required substantially richer capabilities such as explicit planning, tool use, memory, and reflective reasoning (Ruan et al., 2023). Only recently has a more unified and structured definition begun to emerge (e.g., LLM-based agent = LLM + reasoning + planning + memory + tool use + self-improvement + multi-turn interaction + perception, as discussed by Zhang et al. (2025f)), though even this formulation is not universally applicable. Against this historical backdrop, early systems such as MemoryBank (Zhong et al., 2024) and MemGPT (Packer et al., 2023a) framed their contributions as providing *LLM memory*. Yet what they fundamentally addressed were classical agentic challenges, for example enabling an LLM-based conversational agent to track user preferences, maintain dialogue-state information, and accumulate experience across multi-turn interactions. Under a modern and more mature understanding of agency, such systems are naturally categorized as instances of *agent memory*.

That said, the subsumption is not absolute. A distinct line of research genuinely concerns *LLM-internal memory*: managing the transformer's key–value (KV) cache, designing long-context processing mechanisms, or modifying model architectures (e.g., RWKV (Peng et al., 2023), Mamba (Gu and Dao, 2024; Lieber et al., 2024), diffusion-based LMs (Nie et al., 2025)) to better retain information as sequence length grows. These works focus on intrinsic model dynamics and typically address tasks that do not require agentic behavior, and thus should be considered outside the scope of agent memory.

**Overlap.** Within our taxonomy, the majority of what has historically been called "LLM memory" corresponds to forms of agent memory. Techniques such as *few-shot prompting* (Prabhumoye et al., 2022; Ma et al., 2023a) can be viewed as a form of long-term memory, where past exemplars or distilled task summaries serve as reusable knowledge incorporated through retrieval or context injection. *Self-reflection* and iterative refinement

methods (Madaan et al., 2023; Mousavi et al., 2023; Han et al., 2025c) naturally align with short-term, inside-trial memory, as the agent repeatedly leverages intermediate reasoning traces or outcomes from prior attempts within the same task. Even *KV compression* and context-window management (Yoon et al., 2024; Jiang et al., 2023), when used to preserve salient information across the course of a single task, function as short-term memory mechanisms in an agentic sense. These techniques all support the agent's ability to accumulate, transform, and reuse information throughout a task's execution.

**Distinctions.** In contrast, memory mechanisms that intervene directly in the model's internal state—such as architectural modifications for longer effective context, cache rewriting strategies, recurrent-state persistence, attention-sparsity mechanisms, or externalized KV-store expansions—are more appropriately classified as *LLM memory* rather than agent memory. Their goal is to expand or reorganize the representational capacity of the underlying model, not to furnish a decision-making agent with an evolving external memory base. They do not typically support cross-task persistence, environment-driven adaptation, or deliberate memory operations (e.g., formation, evolution, retrieval), and therefore lie outside the operational scope of agent memory as defined in this survey.

### 2.3.2 Agent Memory vs. RAG

At a conceptual level, *agent memory* and *retrieval-augmented generation* (RAG) exhibit substantial overlap: both systems construct, organize, and leverage auxiliary information stores to extend the capabilities of LLM/agents beyond their native parametric knowledge. For instance, structured representations such as knowledge graphs and indexing strategies appear in both communities' methods, and recent developments in agentic RAG demonstrate how autonomous retrieval mechanisms can interact with dynamic databases in ways reminiscent of agent memory architectures (Singh et al., 2025). Indeed, the engineering stacks underlying many RAG and agent memory systems share common building blocks, including vector indices, semantic search, and context expansion modules.

Despite these technological convergences, the two paradigms have *historically* been distinguished by the contexts in which they are applied. Classical RAG techniques primarily augment an LLM with access to **static knowledge sources**, whether flat document stores, structured knowledge bases, or large corpora externally indexed to support retrieval on demand (Zhang et al., 2025p; Han et al., 2025b). These systems are designed to ground generation in up-to-date facts, mitigate hallucinations, and improve accuracy in knowledge-intensive tasks, but they generally do not maintain an internal, evolving memory of past interactions. In contrast, agent memory systems are instantiated within an agent's **ongoing interaction with an environment**, continuously incorporating new information generated by the agent's own actions and environmental feedback into a persistent memory base (Wang et al., 2024l; Zhao et al., 2024; Sun et al., 2025d).

In early formulations the distinction between RAG and agent memory was relatively clear: RAG retrieved from externally maintained knowledge for a single task invocation, whereas agent memory evolved over multi-turn, multi-task interaction. However, this boundary has become increasingly blurred as retrieval systems themselves become more dynamic. For example, certain retrieval tasks continuously update relevant context during iterative querying (e.g., multi-hop QA settings where related context is progressively added). Interestingly, systems such as HippoRAG/HippoRAG2 (Gutierrez et al., 2024; Gutiérrez et al., 2025) have been interpreted by both RAG and memory communities as addressing long-term memory challenges for LLMs. Consequently, a more practical (though not perfectly separable) distinction lies in the **task domain**. RAG is predominantly applied to augment LLMs with large, externally sourced context for individual inference tasks, exemplified by classical multi-hop and knowledge-intensive benchmarks such as HotpotQA (Yang et al., 2018), 2WikiMQA (Ho et al., 2020), and MuSiQue (Trivedi et al., 2022). By contrast, agent memory systems are typically evaluated in settings requiring sustained multi-turn interaction, temporal dependency, or environment-driven adaptation. Representative benchmarks include long-context dialogue evaluations such as LoCoMo (Maharana et al., 2024) and LongMemEval (Wu et al., 2025a), complex problem-solving and deep-research benchmarks such as GAIA (Mialon et al., 2023), XBench (Chen et al., 2025b), and BrowseComp (Wei et al., 2025b), code-centric agentic tasks such as SWE-bench Verified (Jimenez et al., 2024), as well as lifelong learning benchmarks such as StreamBench (Wu et al., 2024a). We provide a comprehensive summary of memory-related benchmarks in Section 6.1.

Nevertheless, even this domain-based distinction contains substantial gray areas. Many works self-described as agent memory systems are evaluated under long-document question-answering tasks such as HotpotQA (Wang et al., 2025g,o), while numerous papers foregrounded as RAG systems in fact implement forms of agentic self-improvement, continually distilling and refining knowledge or skills over time. As a result, titles, methodologies, and empirical evaluations frequently blur the conceptual boundary between the two paradigms. To further clarify these relationships, the following three paragraphs draw upon established taxonomies of RAG from (Mei et al., 2025): *modular RAG*, *graph RAG*, and *agentic RAG*, and examine how the core techniques associated with each lineage manifest within both RAG and agent memory systems.

**Modular RAG** Modular RAG refers to architectures in which the retrieval pipeline is decomposed into clearly specified components, such as indexing, candidate retrieval, reranking, filtering, and context assembly, that operate in a largely static and pipeline-like fashion (Singh et al., 2025). These systems treat retrieval as a well-engineered, modular subsystem external to the LLM, designed primarily for injecting relevant knowledge into the model's context window during inference. Within the agent memory perspective, the corresponding techniques typically appear in the *retrieval stage*, where memory access is realized through vector search, semantic similarity matching, or rule-based filtering, as seen in popular agent memory frameworks like Memary (Memary, 2025), MemOS (Li et al., 2025k), and Mem0 (Chhikara et al., 2025).

**Graph RAG** Graph RAG systems structure the knowledge base as a graph, ranging from knowledge graphs to concept graphs or document-entity relations, and leverage graph traversal or graph-based ranking algorithms to retrieve context (Peng et al., 2024). This representation enables multi-hop relational reasoning, which has proven effective for knowledge-intensive tasks (Edge et al., 2025; Han et al., 2025b; Dong et al., 2025a). In the context of agent memory, graph-structured memory arises naturally when agents accumulate relational insights over time, such as linking concepts, tracking dependencies among subtasks, or recording causal relations inferred through interaction. Several well-established practices include Mem0$^g$ (Chhikara et al., 2025), A-MEM (Xu et al., 2025c), Zep (Rasmussen et al., 2025), and G-memory (Zhang et al., 2025c). Notably, graph-based agent memory systems may *construct, extend, or reorganize* its internal graph throughout the agent's operation. Consequently, graph-based retrieval forms the structural backbone for both paradigms, but only agent memory treats the graph as a living, evolving representation of experience. We provide further analysis on graph-based memory forms in Section 3.1.2 and also refer the readers to a relevant survey (Liu et al., 2025g).

**Agentic RAG** Agentic RAG integrates retrieval into an autonomous decision-making loop, where an LLM agent actively controls when, how, and what to retrieve (Singh et al., 2025; Sun et al., 2025d). These systems often employ iterative querying, multi-step planning, or self-directed search procedures, enabling the agent to refine its information needs through deliberate reasoning, as implemented in PlanRAG (Lee et al., 2024b) and Self-RAG (Asai et al., 2023). For a more detailed understanding of agentic RAG, we refer the readers to Singh et al. (2025). From the agent memory perspective, agentic RAG occupies the closest conceptual space: both systems involve autonomous interaction with an external information store, both support multi-step refinement, and both may incorporate retrieved insights into subsequent reasoning. The key distinction is that classical agentic RAG typically operates over an *external* and often task-specific database, whereas agent memory maintains an *internal, persistent, and self-evolving* memory base that accumulates knowledge across tasks (Yan et al., 2025a; Xu et al., 2025c).

### 2.3.3 Agent Memory vs. Context Engineering

The relationship between *agent memory* and *context engineering* is best understood as an intersection of distinct operational paradigms rather than a hierarchical subsumption. Context engineering is a systematic design methodology that treats the context window as a constrained computational resource. It rigorously optimizes the information payload, including instructions, knowledge, state, and memory, to mitigate the asymmetry between massive input capacity and the model's generation capability (Mei et al., 2025). While agent memory focuses on the **cognitive modeling** of a persistent entity with an evolving identity, context engineering operates under a **resource management** paradigm. From the perspective of context engineering, agent memory is merely one variable within the context assembly function that requires efficient scheduling to maximize inference efficacy. Conversely, from the perspective of an agent, context engineering serves as the

implementation layer that ensures cognitive continuity remains within the physical limits of the underlying model.

**Overlap**   The two fields converge significantly in the technical realization of working memory during long-horizon interactions and often employ functionally identical mechanisms to address the constraints imposed by a finite context window (Hu et al., 2025a; Zhang et al., 2025q; Kang et al., 2025c; Yu et al., 2025a). Both paradigms rely on advanced information compression (Zhou et al., 2025b; Wu et al., 2025f), organization (Xu et al., 2025c; Zhang et al., 2025c; Anokhin et al., 2024), and selection (Zhang et al., 2025q) techniques to preserve operational continuity over extended interaction sequences. For example, token pruning and importance-based selection methods (Jiang et al., 2023; Li et al., 2023c) that are central to context engineering frameworks play a fundamental role in agentic memory systems by filtering noise and retaining salient information. Similarly, the rolling summary technique serves as a shared foundational primitive, functioning simultaneously as a buffer management strategy and a transient episodic memory mechanism (Yu et al., 2025a; Lu et al., 2025b). In practice, the boundary between engineering the context and maintaining an agent's short-term memory effectively dissolves in these scenarios, as both rely on the same underlying summarization, dynamic information retrieval, and recursive state updates (Tang et al., 2025b; Yoon et al., 2024).

**Distinctions**   The distinction becomes most pronounced when moving beyond short-term text processing to the broader scope of long-lived agents. Context engineering primarily addresses the *structural organization* of the interaction interface between LLMs and their operational environment. This includes optimizing tool-integrated reasoning and selection pipelines (Qin et al., 2024a; Schick et al., 2023; Jia and Li, 2025) and standardizing communication protocols, such as MCP (Qiu et al., 2025c). These methods focus on ensuring that instructions, tool calls, and intermediate states are correctly formatted, efficiently scheduled, and executable within the constraints of the context window. As such, context engineering operates at the level of *resource allocation and interface correctness*, emphasizing syntactic validity and execution efficiency.

In contrast, agent memory defines a substantially broader cognitive scope. Beyond transient context assembly, it encompasses the persistent storage of factual knowledge (Zhong et al., 2024), the accumulation and evolution of experiential traces (Zhao et al., 2024; Tang et al., 2025d; Zhang et al., 2025d), and, in some cases, the internalization of memory into model parameters (Wang et al., 2025n). Rather than managing how information is presented to the model at inference time, agent memory governs what the agent *knows*, what it *has experienced*, and how these elements evolve over time. This includes consolidating repeated interactions into knowledge (Tan et al., 2025c), abstracting procedural knowledge from past successes and failures (Ouyang et al., 2025), and maintaining a coherent identity across tasks and episodes (Wang et al., 2024f).

From this perspective, context engineering constructs the external scaffolding that enables perception and action under resource constraints, whereas agent memory constitutes the internal substrate that supports learning, adaptation, and autonomy. The former optimizes the momentary interface between the agent and the model, while the latter sustains a persistent cognitive state that extends beyond any single context window.

## 3   Form: What Carries Memory?

As a starting point for organizing prior work, we begin by examining the most fundamental representational units out of which agent memory can be constructed. We first try to answer: what architectural or representational forms can agent memory take?

Across diverse agent systems, memory is not realized through a single, unified structure. Instead, different task settings call for different storage forms, each with its own structural properties. These architectures endow memory with distinct capabilities, shaping how an agent accumulates information over interactions and maintains behavioral consistency. They ultimately enable memory to fulfill its intended roles across varied task scenarios.

Based on where memory resides and in what form it is represented, we organize these memories into three categories:

> **Three Major Memory Forms**
>
> 1. **Token-level Memory** (Section 3.1): Memory organized as explicit and discrete units that can be individually accessed, modified, and reconstructed. These units remain externally visible and can be stored in a structured form over time.
>
> 2. **Parametric Memory** (Section 3.2): Memory stored within the model parameters, where information is encoded through the statistical patterns of the parameter space and accessed implicitly during forward computation.
>
> 3. **Latent Memory** (Section 3.3): Memory represented in the model's internal hidden states, continuous representations, or evolving latent structures. It can persist and update during inference or across interaction cycles, capturing context-dependent internal states.

The three memory forms outlined above establish the core structural framework for understanding "what carries memory". Each form organizes, stores, and updates information in its own way, giving rise to distinct representational patterns and operational behaviors. With this structural taxonomy in place, we can more systematically examine why agents need memory (Section 4) and how memory evolves, adapts, and shapes agent behavior over sustained interactions (Section 5). This classification provides the conceptual foundation for the discussions that follow.

## 3.1 Token-level Memory

> **Definition of Token-level Memory**
>
> Token-level memory stores information as persistent, discrete units that are externally accessible and inspectable. The token here is a broad representational notion: beyond text tokens, it includes visual tokens, audio frames—any discrete element that can be written, retrieved, reorganized, and revised outside model parameters.

Because these units are explicit, token-level memory is typically transparent, easy to edit, and straightforward to interpret, making it a natural layer for retrieval, routing, conflict handling, and coordination with parametric and latent memory. Token-level memory is also the most common memory form and the one with the largest body of existing work.

Although all token-level memories share the property of being stored as discrete units, they differ significantly in how these units are organized. The structural organization of stored tokens plays a central role in determining how efficiently the agent can search, update, or reason over past information. To describe these differences, we categorize token-level memory by inter-unit structural organization, moving from no explicit topology to multi-layer topologies:

> **Three Major Types of Token-level Memory**
>
> 1. **Flat Memory (1D)**: No explicit inter-unit topology. Memories are accumulated as sequences or bags of units (e.g., snippets, trajectories, chunks)
>
> 2. **Planar Memory (2D)**: A structured but single-layer organization within one plane: units are related by a graph, tree, table and so on, with no cross-layer relations. The structure is explicit, but not layered.
>
> 3. **Hierarchical Memory (3D)**: Structured across multiple layers with inter-layer links, forming a volumetric or stratified memory

The three types of token-level memory are clearly illustrated in Figure 3. From Flat Memory with no topology, to Planar Memory with single-layer structural organization, to Hierarchical Memory with multi-layer interlinked structures, this organizational spectrum governs not only how token-level memory supports search,

**Figure 3** Taxonomy of token-level memory organized by topological complexity and dimensionality: (a) **Flat Memory (1D)** stores information as linear sequences or independent clusters without explicit inter-unit topology, commonly used for *Chunk* sets, *Dialogue* logs, and *Experience* pools. (b) **Planar Memory (2D)** introduces a single-layer structured layout where units are linked via **Tree** or **Graph** structures to capture relational dependencies, supporting diverse node types such as images and chat records. (c) **Hierarchical Memory (3D)** employs multi-level forms, such as **Pyramids** or **Multi-layer** graphs, to facilitate vertical abstraction and cross-layer reasoning between different data granularities, such as raw docs and synthesized QAs.

update, and reasoning, but also how the memory itself is structured and what capabilities it affords. In the subsections that follow, we introduce each organizational form in terms of its strengths and limitations, typical use cases, and representative work. The summary and comparison of representative token-level memory methods are presented in Table 1.

It is worth noting that, following the idea introduced by ReAct (Yao et al., 2023b), a series of studies began focusing on long-horizon interaction tasks (Song et al., 2025a; Jin et al., 2025; Li et al., 2025g,e,i; Wu et al., 2025b). Many of these tasks introduce an explicit notion of memory, and because the memory is generally stored in plaintext form, they fall within the scope of token-level memory. Most of them emphasize how to compress or fold accumulated interaction traces so that agents can operate over long sequences without exceeding context limits (Zhou et al., 2025b; Zhang et al., 2025q; Wu et al., 2025f; Sun et al., 2025a; Li et al., 2025h; Chen et al., 2025a). A more detailed discussion is provided in Section 4.3 about working memory.

### 3.1.1 Flat Memory (1D)

> **Definition of Flat (1D) Memory**
>
> Flat Memory stores information as accumulations of discrete units, without explicitly modeling semantic or relational dependencies among them. These units may include text chunks, user profiles, experience trajectories, their corresponding vector representations, or multimodal entries. Relationships among these units are not encoded directly in the memory.

To facilitate a clear and coherent presentation, we group prior work on flat memory according to their primary design objectives and technical emphases. This grouping serves **an organizational purpose** and does not imply

that the resulting categories are strictly parallel or mutually exclusive. In practice, certain methods may be applicable to multiple categories, and some approaches involving multimodal information may be discussed in other sections when multimodality is not their central focus. Such an organization allows us to systematically review the literature while preserving flexibility in interpretation.

**Table 1** Comparison of representative token-level memory methods. We categorize existing works into three groups based on their topological complexity: **Flat Memory (1D)** for linear or independent records, **Planar Memory (2D)** for structured single-layer graphs/trees, and **Hierarchical Memory (3D)** for multi-level architectures. Methods are characterized across four dimensions: (1) **Multi** indicates multimodal capability, where ✔ denotes support for modalities beyond text (e.g., visual) and ✗ implies text-only; (2) **Type** identifies the specific functional category of the memory (e.g., *Fact* for factual memory, *Exp* for experiential memory, *Work* for working memory ); (3) **Memory Form** details the content of the stored units; and (4) **Task** lists the primary application domains.

| Method | Multi | Type | Memory Form | Task |
|---|---|---|---|---|
| *Flat Memory Models* | | | | |
| Reflexion (Shinn et al., 2023b) | ✗ | E&W | Trajectory as short-term and feedback as long-term | QA, Reasoning, Coding |
| Memento (Zhou et al., 2025a) | ✗ | Exp | Trajectory case (success/failure). | Reasoning |
| JARVIS-1 (Wang et al., 2025p) | ✔ | Exp | Plan-environment pairs. | Game |
| Expel (Zhao et al., 2024) | ✗ | Exp | Insights and few-shot examples. | Reasoning |
| Buffer of Thoughts (Yang et al., 2024b) | ✗ | Exp | High-level thought-templates. | Game, Reasoning, Coding |
| SAGE (Liang et al., 2025) | ✗ | Exp | Dual-store with forgetting mechanism. | Game, Reasoning, Coding |
| ChemAgent (Tang et al., 2025c) | ✗ | Exp | Structured sub-tasks and principles. | Chemistry |
| AgentKB (Tang et al., 2025d) | ✗ | Exp | 5-tuple experience nodes. | Coding, Reasoning |
| H$^2$R (Ye et al., 2025b) | ✗ | Exp | Planning and Execution layers. | Game, Embodied Simulation |
| AWM (Wang et al., 2024l) | ✗ | Exp | Abstracted universal workflows. | Web |
| PRINCIPLES (Kim et al., 2025a) | ✗ | Exp | Rule templates from self-play. | Emotional Companion |
| ReasoningBank (Ouyang et al., 2025) | ✗ | Exp | Transferable reasoning strategy items. | Web |
| Voyager (Wang et al., 2024b) | ✔ | Exp | Executable skill code library. | Game |
| DGM (Zhang et al., 2025h) | ✗ | Exp | Recursive self-modifiable codebase. | Coding |
| Memp (Fang et al., 2025d) | ✗ | Exp | Instructions and abstract scripts. | Embodied Simulation, Travel Planning |
| UFO2 (Zhang et al., 2025a) | ✔ | Exp | System docs and interaction records. | Windows OS |
| LEGOMem (Han et al., 2025a) | ✗ | Exp | Vectorized task trajectories. | Office |
| ToolMem (Xiao et al., 2025b) | ✗ | Exp | Tool capability. | Tool Calling |
| SCM (Wang et al., 2025a) | ✗ | Fact | Memory stream and vector database. | Long-context |
| MemoryBank (Zhong et al., 2024) | ✗ | Fact | History and user profile. | Emotional Companion |
| MPC (Lee et al., 2023) | ✗ | Fact | Persona and summary vector pool. | QA |
| RecMind (Wang et al., 2024h) | ✗ | Fact | User metadata and external knowledge. | Recommendation |
| InteRecAgent (Huang et al., 2025d) | ✗ | Fact | User profiles and candidate item. | Recommendation |
| Ego-LLaVA (Shen et al., 2024) | ✔ | Fact | Language-encoded chunk embeddings. | Multimodal QA |
| ChatHaruhi (Li et al., 2023a) | ✗ | Fact | Dialogue database from media. | Role-Playing |
| Memochat (Lu et al., 2023) | ✗ | Fact | Memos and categorized dialogue history. | Long-conv QA |
| RecursiveSum (Wang et al., 2025h) | ✗ | Fact | Recursive summaries of short dialogues. | Long-conv QA |
| MemGPT (Packer et al., 2023a) | ✗ | Fact | Virtual memory (Main/External contexts). | Long-conv QA, Doc QA |
| RoleLLM (Wang et al., 2024d) | ✗ | Fact | Role-specific QA pairs. | Role-Playing |
| Think-in-memory (Liu et al., 2023a) | ✗ | Fact | Hash table of inductive thoughts. | Long-conv QA |
| PLA (Yuan et al., 2025b) | ✗ | Fact | Evolving records of history and summaries. | QA, Human Feedback |
| COMEDY (Chen et al., 2025c) | ✗ | Fact | Single-model compressed memory format. | Summary, Compression, QA |
| Memoro (Zulfikar et al., 2024) | ✔ | Fact | Speech-to-text vector embeddings. | User Study |
| Memory Sharing (Gao and Zhang, 2024a) | ✗ | Fact | Query-Response pair retrieval. | Literary Creation, Logic, Plan Generation |
| Conv Agent(Alonso et al., 2024) | ✗ | Fact | Chain-of-tables and vector entries. | QA |
| EM-LLM (Fountas et al., 2025) | ✗ | Fact | Episodic events with Bayesian boundaries. | Long-context |
| Memocrs (Xi et al., 2024a) | ✗ | Fact | User metadata and knowledge. | Recommendation |
| SECOM (Pan et al., 2025) | ✗ | Fact | Paragraph-level segmented blocks. | Long-conv QA |
| Mem0 (Chhikara et al., 2025) | ✗ | Fact | Summary and original dialogue. | Long-conv QA |

**Table 1** Comparison of representative token-level memory methods. We categorize existing works into three groups based on their topological complexity: **Flat Memory (1D)** for linear or independent records, **Planar Memory (2D)** for structured single-layer graphs/trees, and **Hierarchical Memory (3D)** for multi-level architectures. Methods are characterized across four dimensions: (1) **Multi** indicates multimodal capability, where ✔ denotes support for modalities beyond text (e.g., visual) and ✗ implies text-only; (2) **Type** identifies the specific functional category of the memory (e.g., *Fact* for factual memory, *Exp* for experiential memory, *Work* for working memory ); (3) **Memory Structure** details the organization mechanism of the stored units; and (4) **Task** lists the primary application domains. (continued)

| Method | Multi | Type | Memory Structure | Task |
|---|---|---|---|---|
| RMM (Tan et al., 2025c) | ✗ | Fact | Reflection-organized flat entries. | Personalization |
| MEMENTO (Kwon et al., 2025) | ✔ | Fact | Interaction history entries. | Personalization |
| MemGuide (Du et al., 2025b) | ✗ | Fact | Dialogue-derived QA pairs. | Long-conv QA |
| MIRIX (Wang and Chen, 2025) | ✔ | Fact | Six optimized flat memory types. | Long-conv QA |
| SemanticAnchor (Chatterjee and Agarwal, 2025) | ✗ | Fact | Syntactic 5-tuple structure. | Long-conv QA |
| MMS (Zhang et al., 2025b) | ✗ | Fact | Dual Retrieval and Context units. | Long-conv QA |
| Memory-R1 (Yan et al., 2025b) | ✗ | Fact | RL-managed mem0 architecture. | Long-conv QA |
| ComoRAG (Wang et al., 2025f) | ✗ | Fact | Fact/Semantic/Plot units with probes. | Narrative QA |
| Nemori (Nan et al., 2025) | ✗ | Fact | Predictive calibration store. | Long-conv QA |
| Livia (Xi and Wang, 2025) | ✔ | Fact | Pruned interaction history. | Emotional Companion |
| MOOM (Chen et al., 2025d) | ✗ | Fact | Decoupled plot and character stores. | Role-Playing |
| Mem-$\alpha$ (Wang et al., 2025o) | ✗ | Fact | Core, Semantic, and Episodic Mem. | Memory Management |
| Personalized Long term Interaction (Westhäußer et al., 2025) | ✗ | Fact | Hierarchical history and summaries. | Personalization |
| LightMem (Fang et al., 2025b) | ✗ | Fact | Optimized Long/Short-term store. | Long-conv QA |
| MEXTRA (Wang et al., 2025b) | ✗ | Fact | Extracted raw dialogue data. | Privacy Attack |
| MovieChat (Song et al., 2024) | ✔ | Fact | Short-term features and long-term persistence. | Video Understanding |
| MA-LMM (He et al., 2024) | ✔ | Fact | Visual and Query memory banks. | Video Understanding |
| VideoAgent (Wang et al., 2024g) | ✔ | Fact | Temporal text descriptions and object tracking. | Video Understanding |
| KARMA (Wang et al., 2025q) | ✔ | Fact | 3D scene graph and dynamic object states. | Embodied Task |
| Embodied VideoAgent (Fan et al., 2025) | ✔ | Fact | Persistent object and sensor store. | MultiModal |
| Mem2Ego (Zhang et al., 2025l) | ✔ | Fact | Map, landmark, and visited location stores. | Embodied Navigation |
| Context-as-Memory (Yu et al., 2025b) | ✔ | Fact | Generated context frames. | Video Generation |
| RCR-Router (Liu et al., 2025c) | ✗ | Fact | Budget-aware semantic subsets. | QA |
| *Planar Memory Models* | | | | |
| D-SMART (Lei et al., 2025) | ✗ | Fact | Structured memory with reasoning trees. | Long-conv QA |
| Reflexion (Shinn et al., 2023b) | ✗ | Work | Reflective text buffer from experiences. | QA, Reasoning, Coding |
| PREMem (Kim et al., 2025b) | ✗ | Fact | Dynamic cross-session linked triples. | Long-conv QA |
| Query Reconstruct (Xu et al., 2025b) | ✗ | Exp | Logic graphs built from knowledge bases. | KnowledgeGraph QA |
| KGT (Sun et al., 2024) | ✗ | Fact | KG node from query and feedback. | QA |
| Optimus-1 (Li et al., 2024d) | ✔ | F&E | Knowledge graph and experience pool. | Game |
| SALI (Pan et al., 2024) | ✔ | Exp | Topological graph with spatial nodes | Navigation |
| HAT (A et al., 2024) | ✗ | Fact | Hierarchical aggregate tree. | Long-conv QA |
| MemTree (Rezazadeh et al., 2025c) | ✗ | Fact | Dynamic hierarchical conversation tree. | Long-conv QA |
| TeaFarm (iunn Ong et al., 2025) | ✗ | Fact | Causal edges connecting memories. | Long-conv QA |
| COMET (Kim et al., 2024b) | ✗ | Fact | Context-aware memory through graph. | Long-conv QA |
| Intrinsic Memory (Yuen et al., 2025) | ✗ | Fact | Private internal and shared external mem. | Planning |
| A-MEM (Xu et al., 2025c) | ✗ | Fact | Card-based connected mem. | Long-conv QA |
| Ret-LLM (Modarressi et al., 2023) | ✗ | Fact | Triplet table and LSH vectors. | QA |
| HuaTuo (Wang et al., 2023a) | ✗ | Fact | Medical Knowledge Graph. | Medical QA |
| M3-Agent (Long et al., 2025) | ✔ | Fact | Multimodal nodes in graph structure. | Embodied QA |
| *Hierarchical Memory Models* | | | | |
| GraphRAG (Edge et al., 2025) | ✗ | Fact | Multi-level community graph indices. | QA, Summarization |
| H-Mem (Sun and Zeng, 2025) | ✗ | Fact | Decoupled index layers and content layers. | Long-conv QA |
| EMG-RAG (Wang et al., 2024k) | ✗ | Fact | Three-tiered memory graph. | QA |

**Table 1** Comparison of representative token-level memory methods. We categorize existing works into three groups based on their topological complexity: **Flat Memory (1D)** for linear or independent records, **Planar Memory (2D)** for structured single-layer graphs/trees, and **Hierarchical Memory (3D)** for multi-level architectures. Methods are characterized across four dimensions: (1) **Multi** indicates multimodal capability, where ✔ denotes support for modalities beyond text (e.g., visual) and ✗ implies text-only; (2) **Type** identifies the specific functional category of the memory (e.g., *Fact* for factual memory, *Exp* for experiential memory, *Work* for working memory ); (3) **Memory Structure** details the organization mechanism of the stored units; and (4) **Task** lists the primary application domains. (continued)

| Method | Multi | Type | Memory Structure | Task |
|---|---|---|---|---|
| G-Memory (Zhang et al., 2025c) | ✗ | Exp | Query-centric three-layer graph structure. | QA, Game, Embodied Task |
| Zep (Rasmussen et al., 2025) | ✗ | Fact | Temporal Knowledge Graphs. | Long-conv QA |
| SGMem (Wu et al., 2025h) | ✗ | Fact | Chunk Graph and Sentence Graph. | Long-conv QA |
| HippoRAG (Gutierrez et al., 2024) | ✗ | Fact | Knowledge with query nodes. | QA |
| HippoRAG 2 (Gutiérrez et al., 2025) | ✗ | Fact | KG with phrase and passage. | QA |
| AriGraph (Anokhin et al., 2024) | ✗ | Fact | Semantic and Episodic memory graph. | Game |
| Lyfe Agents (Kaiya et al., 2023) | ✗ | Fact | Working, Short & Long-term layers. | Social Simulation |
| CAM (Li et al., 2025f) | ✗ | Fact | Multilayer graph with topic. | Doc QA |
| HiAgent (Hu et al., 2025a) | ✗ | E&W | Goal graphs with recursive cluster. | Agentic Tasks |
| ILM-TR (Tang et al., 2024) | ✗ | Fact | Hierarchical Memory tree. | Long-context |

**Dialogue**    Some flat memory work focuses on storing and managing dialogue content. Early approaches primarily focused on preventing forgetting by storing raw dialogue history or generating recursive summaries to extend context windows (Wang et al., 2025a; Lu et al., 2023; Wang et al., 2025h; Yuan et al., 2025b). MemGPT (Packer et al., 2023a) introduces an operating-system metaphor with hierarchical management, inspiring subsequent works (Li et al., 2025k; Kang et al., 2025a) to decouple active context from external storage for infinite context management.

To improve retrieval precision, the granularity and structure of memory units have become increasingly diverse and cognitively aligned. Some works, like COMEDY (Chen et al., 2025c), Memory Sharing (Gao and Zhang, 2024a) and MemGuide (Du et al., 2025b) compress information into compact semantic representations or query-response pairs to facilitate direct lookup, while others, like Alonso et al. (2024) and MIRIX (Wang and Chen, 2025) adopt hybrid structures ranging from vector-table combinations to multi-functional memory types. Furthermore, research has begun to define memory boundaries based on cognitive psychology, organizing information through syntactic tuples (Chatterjee and Agarwal, 2025) or segmenting events based on Bayesian surprise and paragraph structures (Fountas et al., 2025; Pan et al., 2025) , thereby matching human-like cognitive segmentation.

As conversational depth increases, memory evolves to store high-level cognitive processes and narrative complexities. Instead of mere factual records, systems like Think-in-Memory (Liu et al., 2023a) and RMM (Tan et al., 2025c) store inductive thoughts and retrospective reflections to guide future reasoning. In complex scenarios such as role-playing or long narratives, approaches like ComoRAG (Wang et al., 2025f) and MOOM (Chen et al., 2025d) decompose memory into factual, plot-level, and character-level components, ensuring the agent maintains a coherent persona and understanding across extended interactions.

Memory has transitioned from static storage to autonomous and adaptive optimization. Mem0(Chhikara et al., 2025) established standardized operations for memory maintenance, laying the foundation for intelligent control. Recent advances introduce reinforcement learning to optimize memory construction (Yan et al., 2025b; Wang et al., 2025o), while other mechanisms focus on dynamic calibration and efficiency, such as predicting missing information (Nan et al., 2025), managing token budgets across multi-agent systems (Liu et al., 2025c) , and reducing redundancy in long-term storage (Fang et al., 2025b).

**Preference**    Some memory systems focus on modeling a user's evolving tastes, interests, and decision patterns, especially in recommendation scenarios where preference understanding is central. Unlike dialogue-centric memory, which focuses on maintaining conversational coherence, preference memory centers on identifying

a user's tastes and tendencies. Early efforts such as RecMind (Wang et al., 2024h) separate user-specific information from external domain knowledge by storing both factual user attributes and item metadata. InteRecAgent (Huang et al., 2025d) folds memory into the recommendation workflow but focuses more on the current candidate set, keeping user profiles and the active item pool to support context-aware recommendations. MR.Rec (Huang et al., 2025b) builds a memory index archiving the full interaction process, storing raw item information and per-category preference summaries. In conversational settings, Memocrs (Xi et al., 2024a) proposes a more structured design with a user-specific memory tracking entities and user attitudes, and a general memory aggregating cross-user knowledge.

**Profile** A subset of flat memory systems focuses on storing and maintaining stable user profiles, character attributes, or long-term identity information so that agents can behave consistently across turns and tasks. MemoryBank (Zhong et al., 2024) represents one of the earliest frameworks in this direction: it organizes dialogue history and event summaries by timestamp, gradually building a user profile that supports accurate retrieval of identity-relevant information. AI Persona (Wang et al., 2024f) makes the memory system process information not only presented in the dialogue context but also from multi-dimensional human-AI interaction dimensions. MPC (Lee et al., 2023) extends this idea by storing real-time persona information and dialogue summaries in a memory pool, keeping conversation behavior aligned with a consistent persona over long interactions. Westhäußer et al. (2025) proposes a more comprehensive profile-maintenance mechanism, combining long-term and short-term memory with automatically generated summaries after each turn to form a mid-term context, allowing the user profile to evolve continuously through interaction.

In virtual role-playing settings, ChatHaruhi (Li et al., 2023a) extracts dialogue from novels and television scripts, enabling the model to maintain character-consistent behavior by retrieving memory. RoleLLM (Wang et al., 2024d) takes a more structured approach by building question–answer pairs to capture character-specific knowledge.

**Experience** Distinct from the static, general knowledge, experience memory stems from the agent's dynamic accumulation during actual interaction tasks, encompassing specific observations, chains of thought, action trajectories, and environmental feedback. It is important to note that this section just provides a brief overview of experiential memory strictly from the perspective of token-level storage; a more comprehensive analysis and detailed discussion of this domain will be presented in Section 4.2.

The most fundamental form of experience memory involves the direct archival of historical behavioral trajectories. This paradigm enables agents to inform current decision-making by retrieving and reusing past instances, encompassing both successful and failed cases (Zhou et al., 2025a; Wang et al., 2025p).

To address the limited generalizability inherent in raw trajectories, a significant body of research focuses on abstracting specific interactions into higher-level, generalized experiences. As one of the earliest and most influential approaches, Reflexion (Shinn et al., 2023b) distinguishes short-term memory as the trajectory history and long-term memory as the feedback produced by the self-reflection model. Certain studies compress complex interaction histories into universal workflows, rule templates, or high-level "thought-templates" to facilitate cross-problem transfer and reuse (Wang et al., 2024l; Kim et al., 2025a; Yang et al., 2024b). Other works emphasize the structural organization and dynamic maintenance of memory. These approaches ensure that stored insights remain adaptable to novel tasks and are efficiently updated by constructing domain-specific structured knowledge bases, employing hierarchical plan-execute memory architectures, or incorporating human-like forgetting and reflection mechanisms (Tang et al., 2025c,d; Ouyang et al., 2025; Ye et al., 2025b; Zhao et al., 2024; Liang et al., 2025).

In contexts involving programming or specific tool utilization, experience memory evolves into executable skills. Within this paradigm, agents consolidate exploration experiences into code repositories, procedural scripts, or tool-usage entries. Leveraging environmental feedback, these systems iteratively refine code quality or even dynamically modify their underlying logic to achieve self-evolution (Wang et al., 2024a; Yin et al., 2025; Fang et al., 2025d; Xiao et al., 2025b). Furthermore, targeting complex environments such as operating systems, some studies distill successful execution records into reusable exemplars or vectorized representations, thereby facilitating an efficient pipeline from offline construction to online allocation (Zhang et al., 2025a; Han et al., 2025a).

**Multimodal**   Multimodal memory systems store information in the form of discrete token-level units extracted from raw multimodal data, such as images, video frames, audio segments, and text, enabling agents to capture, compress, and retrieve knowledge across channels and over long spans of experience. In wearable and egocentric settings, early work such as Ego-LLaVA (Shen et al., 2024) captures first-person video and converts it into lightweight language descriptions. Memoro (Zulfikar et al., 2024) follows a similar philosophy but uses speech-to-text to form embedding-based memory chunks. Building on this direction, Livia (Xi and Wang, 2025) incorporates long-term user memory into an AR system with emotional awareness, applying forgetting curves and pruning strategies.

For video understanding, the emphasis shifts toward separating transient visual cues from enduring contextual information. MovieChat (Song et al., 2024) adopts a short-term/long-term split, storing recent frame features. MA-LMM (He et al., 2024) pushes this further with a dual-bank design—one storing raw visual features and the other retaining query embeddings. VideoAgent (Wang et al., 2024g) adopts a more semantically organized approach, maintaining a temporal memory of textual clip descriptions alongside object-level memory that tracks entities across frames. In interactive video generation, Context-as-Memory (Yu et al., 2025b) shows that simply storing previously generated frames as memory can also be highly effective.

In embodied scenarios, memory becomes inherently tied to spatial structure and ongoing interaction. KARMA (Wang et al., 2025q) introduces a two-tier memory system: long-term memory stores static objects in a 3D scene graph, while short-term memory tracks object positions and state changes. Embodied VideoAgent (Fan et al., 2025) also builds persistent object memories but fuses them with first-person video and additional embodied sensors. Mem2Ego (Zhang et al., 2025l) extends this idea to navigation by separating global maps, landmark descriptions, and visitation histories into three distinct memory stores. Complementing these task-driven designs, MEMENTO (Kwon et al., 2025) provides an evaluation framework that treats multimodal interaction history as an agent's memory, enabling systematic assessment of how well embodied systems utilize accumulated perceptual experience.

**Discussion**   The primary advantage of Flat Memory is their simplicity and scalability: memory can be appended or pruned with minimal cost, and retrieval methods such as similarity search allow flexible access without requiring predefined structure. This makes them suitable for broad recall, episodic accumulation, and rapidly changing interaction histories. However, the lack of explicit relational organization means that coherence and relevance depend heavily on retrieval quality. As the memory grows, redundancy and noise can accumulate, and the model may retrieve relevant units without understanding how they relate, limiting compositional reasoning, long-horizon planning, and abstraction formation. Thus, topology-free collections excel at broad coverage and lightweight updates, but are constrained in tasks requiring structured inference or stable knowledge organization.

### 3.1.2   Planar Memory (2D)

> **Definition of Planar (2D) Memory**
>
> Planar Memory introduces an explicit organizational topology among memory units, but only within a single structural layer, which for short called *2D*. The topology may be a graph, tree, table, implicit connection structure and so on, where relationships such as adjacency, parent–child ordering, or semantic grouping are encoded within one plane, without hierarchical levels or cross-layer references.

The core of Planar memory forms lies in breaking through a single storage pool by establishing explicit association mechanisms, achieving a leap from mere "storage" to "organization".

**Tree**   Tree structures organize information hierarchically and can handle different levels of abstraction. HAT (A et al., 2024) builds a Hierarchical Aggregate Tree by segmenting long interactions and then aggregating them step by step. This multi-level structure supports coarse-to-fine retrieval and performs better than flat vector indices in long-context question answering. To reduce dialogue fragmentation, MemTree (Rezazadeh et al., 2025c) introduces a dynamic representation that infers hierarchical schemas from isolated conversation

logs. It gradually summarizes concrete events into higher-level concepts, allowing agents to use both detailed memories and abstract knowledge.

**Graph**  Graph structures dominate the landscape of 2D memory due to their ability to capture complex associations, causality, and temporal dynamics. Foundational works like Ret-LLM (Modarressi et al., 2023) abstract external storage into addressable triple-based units, enabling the LLM to interact with a relation-centric table that functions like a lightweight knowledge graph. In the medical domain, HuaTuo (Wang et al., 2023a) injects professional knowledge by integrating a structured corpus of Chinese medical knowledge graphs and clinical texts to fine-tune the base model. KGT (Sun et al., 2024) introduces a real-time personalization mechanism where user preferences and feedback are encoded as nodes and edges in a user-specific knowledge graph. For reasoning-intensive tasks, PREMem (Kim et al., 2025b) shifts part of the inference burden to the memory construction phase, deriving structured memory items and their evolution relations from raw dialogue. Similarly, Memory-augmented Query Reconstruction (Xu et al., 2025b) maintains a dedicated query memory that records past KG queries and reasoning steps, using retrieved records to reconstruct more accurate queries. Building on a timeline perspective, TeaFarm (iunn Ong et al., 2025) organizes dialogue history along segmented timelines and applies structured compression to manage lifelong context. COMET (Kim et al., 2024b) further refines conversational memory by using external commonsense bases to parse dialogue and dynamically update a context-aware persona graph with inferred hidden attributes. A-Mem (Xu et al., 2025c) standardizes knowledge into card-like units. It organizes them by relevance and places related memories in the same box, which builds a complete memory network. Intrinsic Memory Agents (Yuen et al., 2025) employ a partitioned architecture in which sub-agents maintain their own role-specific private memories while collaboratively reading and writing to a shared memory. Extending to multimodel agents, M3-Agent (Long et al., 2025) unifies image, audio, and text into an entity-centric memory graph. SALI (Pan et al., 2024) constructs a Reality–Imagination Hybrid Memory, unifying real observations and imagined future scenarios into a consistent navigation graph.

**Hybrid**  Complex tasks often require hybrid architectures that segregate distinct cognitive functions while sharing a common memory substrate. Optimus-1 (Li et al., 2024d) explicitly separates static knowledge into a hierarchical directed knowledge graph for planning, and dynamic interactions into an abstract multimodal experience Pool for reflection and self-improvement. D-SMART (Lei et al., 2025) combines a structured factual memory, implemented as a continuously updated knowledge graph, with a traversal-based reasoning tree.

**Discussion**  The Planar Memory, by effectively establishing links between its nodes, enables memories to leverage collective synergies and thus encode more comprehensive contextual knowledge. Moreover, it supports retrieval mechanisms that go beyond simple iteration, including structured key–value lookups and relational traversal along graph edges. These capabilities make the form strong in storing, organizing, and managing memories. However, it also faces a critical limitation: Without a hierarchical storage mechanism, all memories must be consolidated into a single, monolithic module. As task scenarios grow in complexity and diversity, this redundant and flattened design becomes increasingly inadequate for robust performance. More importantly, the high construction and search costs significantly hinder its practical deployment.

### 3.1.3  Hierarchical Memory (3D)

> **Definition of Hierarchical (3D) Memory**
>
> Hierarchical memory organizes information across layers, using inter-level connections to shape the memories into a volumetric structured space.

Such hierarchies support representations at different degrees of abstraction—from raw observations, to compact event summaries, to higher-level thematic patterns. Cross-layer connections further yield a volumetric memory space through which the system can navigate not only laterally among units but also vertically across abstraction levels.

Hierarchical Memory moves beyond simple stratification, aiming to build complex systems with deep abstraction capabilities and dynamic evolutionary mechanisms. These works typically employ multi-level graph structures or neuroscience-inspired mechanisms to build a more human-like volumetric memory space, where information is richer and the connections between memory units are clearer and more explicit.

**Pyramid** This category constructs memory as multi-level pyramids, where information is progressively organized into higher layers of abstraction and queried in a coarse-to-fine manner. HiAgent (Hu et al., 2025a) manages long-horizon tasks through a subgoal-centered hierarchical working memory, keeping detailed trajectories for the currently active subgoal while compressing completed subgoals into higher-level summaries that can be selectively retrieved when needed. GraphRAG (Edge et al., 2025) builds a multi-level graph index via community detection, recursively aggregating entity-level subgraphs into community-level summaries. Extending the idea of clustering memory nodes, Zep (Rasmussen et al., 2025) formalizes agent memory as a Temporal Knowledge Graph, and it similarly performs community partitioning. ILM-TR (Tang et al., 2024) employs a tree-structured, pyramidal index coupled with an Inner Loop mechanism, repeatedly querying summaries at different abstraction levels and updating a short-term memory buffer until the retrieved evidence and generated answer stabilize. To ensure controllable personalization, EMG-RAG (Wang et al., 2024k) organizes an Editable Memory Graph into three tiers, where a tree-like type and subclass index (L1, L2) sits above an entity-level memory graph (L3). In multi-agent systems, G-Memory (Zhang et al., 2025c) structures shared experience using a three-tier graph hierarchy of insight, query, and interaction graphs. This design enables query-centric traversal to move vertically between high-level cross-trial insights and compact trajectories of concrete collaborations.

**Multi-Layer** These forms instead emphasize layered specialization, organizing memory into distinct modules or levels that focus on particular information types or functions. Lyfe Agents (Kaiya et al., 2023) separates salient long-term records from low-value transient details, allowing the system to maintain a compact, behaviorally important layer of memories. H-Mem (Sun and Zeng, 2025) explicitly arranges long-term dialogue memory into a multi-level hierarchy ordered by semantic abstraction, where lower layers store fine-grained interaction snippets and higher layers store increasingly compressed summaries. Biologically inspired architectures such as HippoRAG (Gutierrez et al., 2024) factor memory into an associative indexing component, implemented as an open knowledge graph, and an underlying passage store, using the graph layer to orchestrate multi-hop retrieval over stored content. Its successor, HippoRAG 2 (Gutiérrez et al., 2025), extends this design into a non-parametric continual-learning setting, enriching the indexing layer with deeper passage integration and online LLM filtering. AriGraph (Anokhin et al., 2024) separates memory by information type within a unified graph, combining a semantic knowledge-graph world model that encodes environment structure with an event-level component that links concrete observations back to the semantic backbone. Similarly, SGMem (Wu et al., 2025h) adds a sentence-graph memory level on top of raw dialogue, representing histories as sentence-level graphs within chunked units. CAM (Li et al., 2025f) layers the reading process itself by incrementally clustering overlapping semantic graphs into a hierarchical schemata structure.

**Discussion** By placing memory nodes at the intersection of hierarchical and relational dimensions, Hierarchical Memory allows different memories to interact and form multi-dimensional synergies. This design helps the system encode knowledge that is more holistic and more deeply contextualized. The form also supports powerful retrieval: it enables complex, multi-path queries that move through relational networks within each layer and across abstraction levels between layers. This ability allows the system to retrieve task-relevant memories with high precision, leading to strong task performance. However, the structure's complexity and its dense information organization create challenges for both retrieval efficiency and overall effectiveness. In particular, ensuring that all stored memories remain semantically meaningful and designing the optimal three-dimensional layout of the system remain difficult and critical problems.

## 3.2 Parametric Memory

In contrast to token-level memory, which stores information as visible and editable discrete units, parametric memory stores information directly in the model's parameters. In this section, we examine methods that embed memory into learnable parameter spaces, allowing the model to internalize and recall information without referring to external storage.

Based on where the memory is stored relative to the core model parameters, we distinguish two primary forms of parametric memory:

> **Two Major Types of Parametric Memory**
>
> 1. **Internal Parametric Memory**: Memory encoded within the original parameters of the model (e.g., weights, biases). These methods directly adjust the base model to incorporate new knowledge or behavior.
>
> 2. **External Parametric Memory**: Memory stored in additional or auxiliary parameter sets, such as adapters, LoRA modules, or lightweight proxy models. These methods introduce new parameters to carry memory without modifying the original model weights.

This distinction reflects a key design choice: whether memory is fully absorbed into the base model or attached modularly alongside it. In the subsections that follow, for each form we outline the implementation methods, analyze its strengths and limitations, and list representative systems or work. Table 2 provides an overview of representative parametric memory methods.

### 3.2.1 Internal Parametric Memory

Internal parameter memory injects domain knowledge, personalized knowledge, or priors required by downstream tasks into the model. We also regard enhancing the model's long-context capability as injecting a prior. The timing of memory injection can be the pre-training phase, continued pre-training phase, mid-training phase, or post-training phase. The memory stored in internal parameters does not add extra parameters or additional modules.

**Pre-Train**   Some works introduce memory mechanisms during the pre-training phase, aiming to address the issue that long-tail world knowledge is difficult to compress into the limited model parameters. LMLM (Zhao et al., 2025b) and HierMemLM (Pouransari et al., 2025) store the memory for knowledge retrieval in the model during the pre-training phase, while storing the knowledge itself in an external knowledge base. Some works also optimize the computational efficiency of attention to enhance long-window memory capability (Xiao et al., 2024; Qin et al., 2024b,c; Dao, 2024; Shah et al., 2024).

**Mid-Train**   During the continued pre-training phase, some works incorporate generalizable experience from downstream tasks. For instance, Su et al. (2025) and Zhang et al. (2025j) integrate agent experience. Some works improve the long-window performance or efficiency of LLMs during the mid-training phase, enabling the model to maintain more short-term memory with longer windows in memory-aided tasks (Zaheer et al., 2020; Chen et al., 2024a).

**Post-Train**   Other works incorporate memory during the post-training phase to adapt to downstream tasks. Some works enable LLMs to memorize personalized user history or styles. Some works allow LLMs to learn from the successes or failures of past similar task executions. Character-LM (Shao et al., 2023) and CharacterGLM (Zhou et al., 2024a) fine-tunes the LLM into different characteristics. During the post-training phase, SELF-PARAM (Wang et al., 2025n) injects additional knowledge through KL divergence distillation without requiring extra parameters. Room (Kim et al., 2023b) stores knowledge externally while save experience internally. KnowledgeEditor (Cao et al., 2021) modifies internal parameters, aiming to alter only the knowledge that requires editing. MEND (Mitchell et al., 2022) achieves fast knowledge editing by using small networks to modify the gradients of large models. PersonalityEdit (Mao et al., 2024) proposes an LLM personality editing dataset based on personality theories in psychology. APP (Ma et al., 2024) employs multiple training objectives to ensure that adjacent knowledge is minimally disturbed during knowledge editing. DINM (Wang et al., 2024c) proposes a model editing method that enables the model to learn to reject such dangerous requests without affecting its normal functions.

**Discussion**   The advantages of internal parameters lie in their simple structure, which does not add extra inference overhead or deployment costs to the vanilla model. Their drawback is the difficulty in updating

**Table 2** Taxonomy of parametric memory methods. We categorize existing works based on the *storage location* relative to the core model: **Internal Parametric Memory** embeds knowledge directly into the original weights, while **External Parametric Memory** isolates information within auxiliary parameter sets. Based on the training **phase**, we performed a secondary classification of the articles. Methods are compared across three technical dimensions: (1) **Type** defines the nature of the memory, (2) **Task** specifies the target downstream application, and (3) **Optimization** denotes the optimization strategy, such as *SFT*, *FT* (fine-tuning) , and *PE* (prompt engineering).

| Method | Type | Task | Optimization |
|---|---|---|---|
| *I. Internal Parametric Memory* | | | |
| **(a) Pre-Train Phase** | | | |
| TNL (Qin et al., 2024b) | Working | QA, Reasoning | SFT |
| StreamingLLM (Xiao et al., 2024) | Working | QA, Reasoning | SFT |
| LMLM (Zhao et al., 2025b) | Factual | QA, Factual Gen | SFT |
| HierMemLM (Pouransari et al., 2025) | Factual | QA, Language Modeling | SFT |
| Function Token (Zhang et al., 2025n) | Factual | Language Modeling | Pretrain |
| **(b) Mid-Train Phase** | | | |
| Agent-Founder (Su et al., 2025) | Experiential | Tool Calling, Deep Research | SFT |
| Early Experience (Zhang et al., 2025j) | Experiential | Tool Calling, Embodied Simulation, Reasoning, Web | SFT |
| **(c) Post-Train Phase** | | | |
| Character-LM (Shao et al., 2023) | Factual | Role Playing | SFT |
| CharacterGLM (Zhou et al., 2024a) | Factual | Role Playing | SFT |
| SELF-PARAM (Wang et al., 2025n) | Factual | QA, Recommendation | KL Tuning |
| Room (Kim et al., 2023b) | Experiential | Embodied Task | RL |
| KnowledgeEditor (Cao et al., 2021) | Factual | QA, Fact Checking | FT |
| Mend (Mitchell et al., 2022) | Factual | QA, Fact Checking, Model Editing | FT |
| PersonalityEdit Mao et al. (2024) | Factual | QA, Model Editing | FT, PE |
| APP (Ma et al., 2024) | Factual | QA | FT |
| DINM (Wang et al., 2024c) | Experiential | QA, Detoxification | FT |
| AlphaEdit (Fang et al., 2025c) | Factual | QA | FT |
| *II. External Parametric Memory* | | | |
| **(a) Adapter-based Modules** | | | |
| MLP-Memory (Wei et al., 2025d) | Factual | QA, Classification, Textual Entailment | SFT |
| K-Adapter (Wang et al., 2021) | Factual | QA, Entity Typing, Classification | SFT |
| WISE (Wang et al., 2024e) | Factual | QA, Hallucination Detection | SFT |
| ELDER (Li et al., 2025d) | Factual | Model Editing | SFT |
| T-Patcher (Huang et al., 2023) | Factual | QA | FT |
| Lin et al. (2025) | Factual | QA | SFT |
| **(b) Auxiliary LM-based Modules** | | | |
| MAC (Tack et al., 2024) | Factual | QA | SFT |
| Retroformer (Yao et al., 2024a) | Experiential | QA, Web Navigation | RL |

internal parameters: storing new memory requires retraining, which is costly and prone to forgetting old memory. Therefore, internal parameter memory is more suitable for large-scale storage of domain knowledge or task priors, rather than short segments of personalized memory or working memory.

### 3.2.2 External Parametric Memory

Storing memory as tokens outside LLMs leads to insufficient understanding of token-form memory content in the input window by the model. Meanwhile, storing memory in the parameters of LLMs has issues, such as

difficulty in updating and conflicts with pre-trained knowledge. Some works adopt a compromise approach, which **introduces memory through external parameters** without altering the original parameters of LLMs.

**Adapter** A common line of external parametric memory methods relies on modules that are attached to a frozen base model. MLP-Memory (Wei et al., 2025d) integrates RAG knowledge with Transformer decoders through MLP. K-Adapter (Wang et al., 2021) injects new knowledge by training task-specific adapter modules while keeping the original backbone unchanged, enabling continual knowledge expansion without interfering with pre-trained representations. WISE (Wang et al., 2024e) further introduces a dual-parameter memory setup—separating pre-trained knowledge and edited knowledge—and a routing mechanism that dynamically selects which parameter memory to use at inference time, thus mitigating conflicts during lifelong editing. ELDER (Li et al., 2025d) advances this direction by maintaining multiple LoRA modules and learning a routing function that adaptively selects or blends them based on input semantics, improving robustness and scalability in long-term editing scenarios. Collectively, these methods leverage additional parameter subspaces to store and retrieve memory in a modular and reversible manner, avoiding the risks of catastrophic interference associated with directly modifying the core model weights.

**Auxiliary LM** Beyond Adapter-based storage, another line of work adopts a more architecturally decoupled form of external parametric memory, where memory is stored in a separate model or external knowledge module. MAC (Tack et al., 2024) compresses the information from a new document into a compact modulation through an amortization network, and stores it in a memory bank. Retroformer (Yao et al., 2024a) proposes a learning paradigm for memorizing the experiences of successes or failures in past task executions.

**Discussion** This external parametric memory approach provides a balance between adaptability and model stability. Because memory is encoded into additional parameter modules, it can be added, removed, or replaced without interfering with the base model's pre-trained representation space. This supports modular updates, task-specific personalization, and controlled rollback, while avoiding the catastrophic forgetting or global weight distortion that may occur in full model fine-tuning.

However, this approach also comes with limitations. External parameter modules must still integrate with the model's internal representation flow, meaning that their influence is indirect and mediated through the model's attention and computation pathways. As a result, the effectiveness of memory injection depends on how well the external parameters can interface with internal parametric knowledge.

## 3.3 Latent Memory

> **Definition of Latent Memory**
>
> Latent memory refers to memory that is carried implicitly in the model's internal representations (e.g., KV cache, activations, hidden states, latent embeddings), rather than being stored as explicit, human-readable tokens or dedicated parameter sets.

Latent avoids exposing memory in plaintext and introduces practically less inference latency, while potentially offering better performance gains by preserving fine-grained contextual signals within the model's own representational space.

As shown in Figure 4, we organize prior work by the origin of latent memory, which means how the latent state is formed and introduced into the agent. We summarize the works in this part in Table 3.

**Figure 4** Overview of Latent Memory integration in LLM agents. Unlike explicit text storage, latent memory operates within the model's internal representational space. The framework is categorized by the origin of the latent state: (a) **Generate**, where auxiliary models synthesize embeddings to interfere with or augment the LLM's forward pass; (b) **Reuse**, which directly propagates prior computational states such as KV caches or intermediate embeddings; and (c) **Transform**, which compresses internal states through token selection, merging, or projection to maintain efficient context.

---

**Three Major Types of Latent Memory**

1. **Generate**: latent memory is produced by an independent model or a module, and then supplied to the agent as reusable internal representations.

2. **Reuse**: latent memory is directly carried over from prior computation, most prominently KV-cache reuse (within or across turns), as well as recurrent or stateful controllers that propagate hidden states.

3. **Transform**: existing latent state is transformed into new representations(e.g., distillation, pooling, or compression), so the agent can retain essentials while reducing latency and context footprint.

---

### 3.3.1 Generate

A major line of work builds memory by **generating new latent representations** rather than reusing or transforming existing activations. In this paradigm, the model or an auxiliary encoder creates compact continuous states. These states may appear as special tokens in the sequence or as standalone vectors. They summarize the essential information from long contexts, task trajectories, or multimodal inputs. The generated latent summaries are then stored, inserted, or used as conditions for later reasoning or decision-making. This enables the system to operate beyond its native context length, maintain task-specific intermediate states, and retain knowledge across episodes without revisiting the original input. Although the concrete forms vary across studies, the underlying idea remains consistent. Memory is explicitly produced through learned encoding or compression, and the resulting latent states serve as reusable memory units that support future inference.

This design choice may also raise potential ambiguity with parametric memory, particularly since many

**Table 3** Taxonomy of latent memory methods. We categorize existing works based on the *origin* of the latent state: **Generate** synthesizes memory via auxiliary modules, **Reuse** propagates internal computational states, and **Transform** compresses, modifies or restructs existing latent state. Methods are compared across three technical dimensions: (1) **Form** specifies the specific data type of the latent memory, (2) **Type** defines the nature of the recorded content (e.g., Working, Factual, and Experiential), and (3) **Task** denotes the target downstream application.

| Method | Form | Type | Task |
|---|---|---|---|
| *I. Generate* | | | |
| **(a) Single Modal** | | | |
| Gist (Mu et al., 2023) | Gist Tokens | Working | Long-context Compression |
| Taking a Deep Breath (Luo et al., 2024) | Sentinel Tokens | Working | Long-context QA |
| SoftCoT (Xu et al., 2025d) | Soft Tokens | Working | Reasoning |
| CARE (Choi et al., 2025) | Memory Tokens | Working | QA, Fact Checking |
| AutoCompressor (Chevalier et al., 2023) | Summary Vectors | Working | QA, Compression |
| MemoRAG (Qian et al., 2025) | Global Semantic States | Working | QA, Summary |
| MemoryLLM (Wang et al., 2024j) | Persistent Tokens | Factual | Long-conv QA, Model Editing |
| M+ (Wang et al., 2025m) | Cross-layer Token Pools | Factual | QA |
| LM2 (Kang et al., 2025b) | Matrix Slots | Working | QA, Reasoning |
| Titans (Behrouz et al., 2025b) | Neural Weights (MLP) | Working | QA, Language Modeling |
| MemGen (Zhang et al., 2025d) | LoRA Fragments | Working, Exp. | QA, Math, Code, Embodied Task, Reasoning |
| EMU (Na et al., 2024) | Embeddings w/ Returns | Factual | Game |
| TokMem (Wu et al., 2025j) | Memory Tokens | Exp. | Funcation calling |
| Nested Learning (Behrouz et al., 2025a) | Nested Optimization | Factual | Language Modeling |
| **(b) Multi-Modal** | | | |
| CoMem (Wu et al., 2025d) | Multimodal Embeddings | Factual | Multimodal QA |
| ACM (Wu et al., 2025e) | Trajectory Embeddings | Working | Web |
| Time-VLM (Zhong et al., 2025) | Patch Embeddings | Working | Video Understanding |
| Mem Augmented RL (Mezghani et al., 2022) | Novelty State Encoder | Working | Visual Navigation |
| MemoryVLA (Shi et al., 2025a) | Perceptual States | Factual, Working | Embodied Task |
| XMem (Cheng and Schwing, 2022) | Key-Value Embeddings | Working | Video Segmentation |
| *II. Reuse* | | | |
| Memorizing Transformers (Wu et al., 2022) | External KV Cache | Working | Language Modeling |
| SirLLM (Yao et al., 2024b) | Entropy-selected KV | Factual | Long-conv QA |
| Memory$^3$ (Yang et al., 2024a) | Critical KV Pairs | Factual | QA |
| FOT (Tworkowski et al., 2023) | Memory-Attention KV | Working | QA, Few-shot learning, Language Modeling |
| LONGMEM (Wang et al., 2023b) | Residual SideNet KV | Working | Language Modeling and Understanding |
| *III. Transform* | | | |
| Scissorhands (Liu et al., 2023b) | Pruned KV | Working | Image classification & generation |
| SnapKV (Li et al., 2024b) | Aggregated Prefix KV | Working | Language Modeling |
| PyramidKV (Cai et al., 2024) | Layer-wise Budget | Working | Language Modeling |
| RazorAttention (Tang et al., 2025a) | Compensated Window | Working | Language Modeling |
| H2O (Zhang et al., 2023) | Heavy Hitter Tokens | Working | QA, Language Modeling |

methods rely on separately trained models to generate latent representations. In this chapter, however, our classification is grounded in the form of memory rather than the learning mechanism. Crucially, although these approaches generate memory through learned encoding, the produced latent representations are explicitly instantiated and reused as independent memory units, rather than being directly embedded into the model's parameters or forward-pass activations. We will return to this distinction when discussing individual methods in detail.

**Single Modal** In the single-modal setting, a major group of methods focuses on long-context processing and language modeling, where models generate a small set of internal representations to replace long raw inputs (Mu et al., 2023; Luo et al., 2024; Xu et al., 2025d; Chevalier et al., 2023; Qian et al., 2025; Wang et al., 2024j, 2025m). A typical strategy is to compress long sequences into a few internal tokens or continuous vectors that can be reused during later inference. For example, Gist (Mu et al., 2023) train a language model to produce a set of gist tokens after processing a long prompt. Luo et al. (2024) introduce a special sentinel token at each chunk boundary and encourage the model to aggregate local semantics into that token. SoftCoT (Xu et al., 2025d) follows a similar direction by generating instance-specific soft tokens from the last hidden state. CARE (Choi et al., 2025) further extends the latent tokens by training a context assessor that compresses retrieved RAG documents into compact memory tokens.

Work such as AutoCompressor (Chevalier et al., 2023) and MemoRAG (Qian et al., 2025) emphasizes vectorized

or standalone latent representations. AutoCompressor (Chevalier et al., 2023) encodes entire long documents into a small number of summary vectors serving as soft prompts, while MemoRAG (Qian et al., 2025) uses an LLM to produce compact hidden-state memories capturing global semantic structure. These approaches not only abstract away from raw text but also transform retrieved or contextualized information into new latent memory units optimized for reuse. To support more persistent memory, MemoryLLM (Wang et al., 2024j) embeds a set of dedicated memory tokens within the model's latent space. M+ (Wang et al., 2025m) extends this idea into a cross-layer long-term memory architecture. LM2 (Kang et al., 2025b) follows a related but structurally distinct direction by introducing matrix-shaped latent memory slots into every layer.

A different branch of work internalizes the generation of latent memory within the model's parameter dynamics. Although these works rely on parameterized modules, their operational memory units remain latent representations, placing them firmly within this category. Titans (Behrouz et al., 2025b) compresses long-range information into an online-updated MLP weight, producing latent vectors during inference. MemGen (Zhang et al., 2025d) dynamically generates latent memory during decoding: two LoRA adapters determine where to insert memory fragments and what latent content to insert. EMU (Na et al., 2024) trains a state encoder to produce latent embeddings annotated with returns and desirability.

**Multi Modal**  In multimodal settings, generative latent memory extends to images, audios and videos, encoding them as compact latent representations. CoMem (Wu et al., 2025d) uses a VLM to compress multimodal knowledge into a set of embeddings that act as plug-and-play memory. Similarly, Wu et al. (2025e) compresses entire GUI interaction trajectories into fixed-length embeddings and injects them into the VLM input space. For temporal modeling, Time-VLM (Zhong et al., 2025) divides video or interaction streams into patches and generates a latent embedding for each patch.

In vision-based navigation, Mezghani et al. (2022) learns a state encoder that maps visual observations into a latent space and constructs an episodic memory containing only novel observations. MemoryVLA (Shi et al., 2025a) maintains a Perceptual–Cognitive Memory Bank that stores both perceptual details and high-level semantics as transformer hidden states. In long-video object segmentation, XMem(Cheng and Schwing, 2022) encodes each frame into key–value latent embeddings and organizes them into a multi-stage memory comprising perceptual, working, and long-term components.

**Discussion**  These single-modal and multimodal approaches share the same fundamental principle: first generate compact latent representations, then maintain and retrieve them as memory entries. The model can actively construct highly information-dense representations tailored to the task, capturing key dynamics, long-range dependencies, or cross-modal relations with minimal storage cost. It also avoids repeatedly processing the full context, enabling more efficient reasoning across extended interactions.

However, the drawbacks are equally evident. The generation process itself may introduce information loss or bias, and the states can drift or accumulate errors over multiple read–write cycles. Moreover, training a dedicated module to generate latent representations introduces additional computational overhead, data requirements, and engineering complexity.

### 3.3.2 Reuse

In contrast to methods that generate new latent representations, another line of work directly **reuses the model's internal activations, primarily the key–value (KV) cache**, as latent memory. These approaches do not transform(modify, compress) the stored KV pairs and instead treat the raw activations from forward passes as reusable memory entries. The main challenge is to determine which KV pairs to keep, how to index them, and how to retrieve them efficiently under long-context or continual-processing demands.

From a cognitive perspective, Gershman et al. (2025) provides conceptual grounding by framing biological memory as a key–value system, where keys function as retrieval addresses and values encode stored content—an abstraction closely aligned with KV-based memory in modern LLMs. Memorizing Transformers (Wu et al., 2022) explicitly store past KV pairs and retrieve them via K-nearest-neighbor search during inference. FOT (Tworkowski et al., 2023) extends this line of work by introducing memory-attention layers that perform KNN-based retrieval over additional KV memories during inference. LONGMEM (Wang et al., 2023b) similarly augments long-range retrieval, employing a lightweight residual SideNet that treats historical KV

embeddings as a persistent memory store. These systems demonstrate how retrieval-aware organization of latent KV states can substantially enhance access to distant information.

**Discussion**  Reuse-type latent memory methods highlight the effectiveness of directly leveraging the model's own internal activations as memory, showing that carefully curated KV representations can serve as a powerful and efficient substrate for long-range retrieval and reasoning.

Their greatest strength lies in preserving the full fidelity of the model's internal activations, ensuring that no information is lost through pruning or compression. This makes them conceptually simple, easy to integrate into existing forms, and highly faithful to the model's original computation. However, raw KV caches grow rapidly with context length, which increases memory consumption and can make retrieval less efficient. The effectiveness of reuse therefore depends heavily on indexing strategies.

### 3.3.3  Transform

Transform-type latent memory methods focus on **modifying, compressing, or restructuring existing latent states** rather than generating entirely new ones or directly reusing raw KV caches. These approaches treat KV caches and hidden activations as malleable memory units, reshaping them through selection, aggregation, or structural transformation. In doing so, they occupy a conceptual middle ground between generate-type and reuse-type memory: the model does not create fresh latent representations, but it also does more than simply replay stored KV pairs.

A major line of work focuses on compressing KV caches while preserving essential semantics. Some methods reduce memory usage by keeping only the most influential tokens. Scissorhands (Liu et al., 2023b) prunes tokens based on attention scores when cache capacity is exceeded, whereas SnapKV (Li et al., 2024b) aggregates high-importance prefix KV representations via a head-wise voting mechanism. PyramidKV (Cai et al., 2024) reallocates KV budgets across layers. SirLLM (Yao et al., 2024b) builds on this perspective by estimating token importance with a token-entropy criterion and selectively retaining only informative KV entries. Memory[3] (Yang et al., 2024a) only stores the most critical attention key–value pairs, significantly shrinking storage requirements. RazorAttention (Tang et al., 2025a) introduces a more explicit compression scheme: it computes the effective attention span of each head, retains only a limited local window, and uses compensation tokens to preserve information from discarded entries. From a more efficiency-oriented perspective, H2O (Zhang et al., 2023) adopts a simpler eviction strategy, retaining only the most recent tokens along with special H2 tokens to reduce memory footprint.
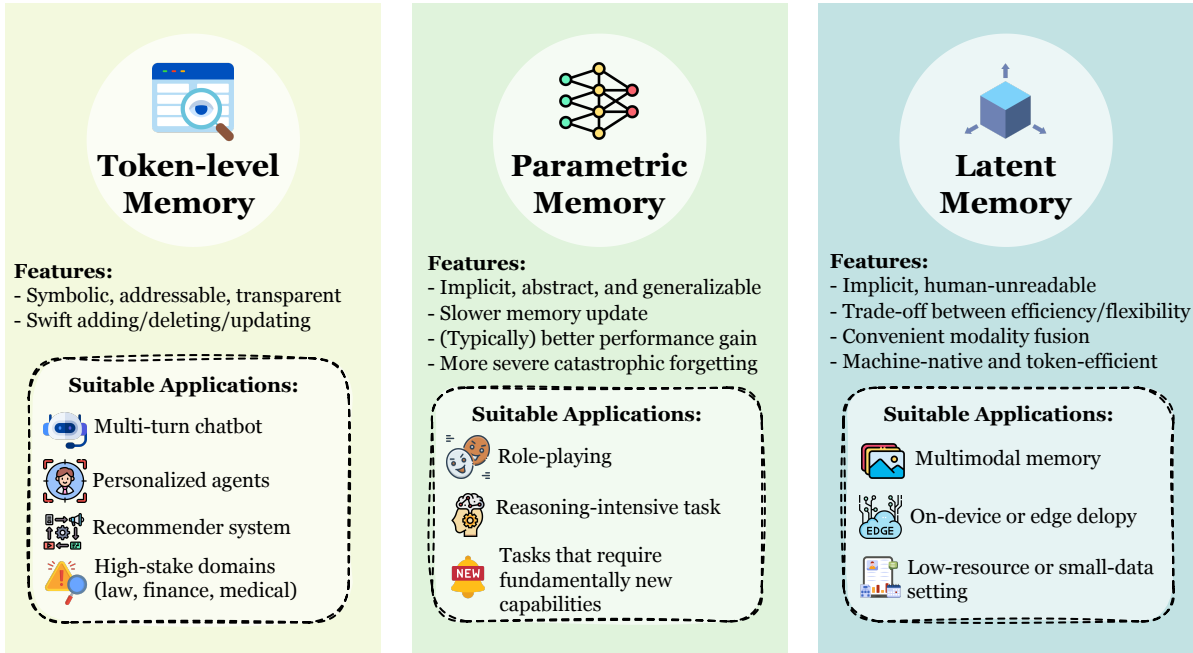
**Discussion**  These methods demonstrate how latent memory can be transformed, through selection, retrieval enhancement, or compressed re-encoding, into more effective memory representations, enabling LLMs to extend their usable context length and improve reasoning performance without relying on raw cache reuse.

Their main advantage lies in producing more compact and information-dense memory representations, which reduce storage cost and enable efficient retrieval over long contexts. By reshaping latent states, these methods allow the model to access distilled semantic signals that may be more useful than raw activations. However, transformation introduces the risk of information loss, and the compressed states can become harder to interpret or verify compared with directly reused KV caches. The additional computation required for pruning, aggregation, or re-encoding also increases system complexity.

## 3.4  Adaptation

As shown above, such a large body of work has focused on agent memory, clearly demonstrating that memory mechanisms are essential for agent systems (Zhang et al., 2025r). The choice of memory type in an agent system reflects how designers expect the agent to behave in a given task. Designers are not simply asking the agent to remember certain information, but also implicitly expressing how they want that information to shape the agent's behavior. Therefore, choosing the right type of memory for a task is far more than a simple combinatorial choice.

In this section, we start from the features of each memory type and discuss which tasks and scenarios they are best suited for in an ideal setting, as shown in Figure 5. We hope this discussion can offer useful ideas

**Figure 5** Overview of three complementary memory paradigms for LLM agents. Token-level, parametric, and latent memories differ in their representational form, update dynamics, interpretability, and efficiency, leading to distinct strengths, limitations, and application domains in long-horizon and interactive agent systems.

and guidance for making practical choices. The examples illustrate only one possible form of memory in these idealized settings and do not imply that other memory types lack unique advantages in the same scenarios.

**Token-level Memory**   Token-level memory remains *symbolic*, *addressable*, and *transparent*, making it particularly well suited for scenarios where explicit reasoning, controllability, and accountability are essential. This type of memory excels in real-time, high-frequency update settings, where an agent must continuously track and revise information, and where the knowledge itself exhibits a clear structure that can be explicitly modeled. Its externalizability allows memory to be easily inspected, audited, transferred, or revised, making it especially suitable for domains requiring precise add/delete/update operations. The high level of interpretability further ensures that an agent's decision process can be traced back to concrete memory units, a crucial property in high-stakes applications. Moreover, token-level memory provides long-term stability and avoids catastrophic forgetting, enabling agents to accumulate reliable knowledge over extended time horizons. Another practical advantage is that token-level memory is often implemented as a plug-and-play module, allowing it to be readily integrated with the latest closed-source or open-source foundation models without modifying their internal parameters.

**Possible Scenarios:**

- Chatbots and multi-turn dialogue systems. (Zhong et al., 2024; Lu et al., 2023; Chhikara et al., 2025)
- Long-horizon or life-long agents requiring stable memory. (Wang et al., 2024f; Westhäußer et al., 2025)
- User-specific personalization profiles. (Wang et al., 2024f; Lee et al., 2023)
- Recommendation systems. (Wang et al., 2024h; Huang et al., 2025d; Xi et al., 2024a)
- Enterprise or organizational knowledge bases.
- Legal, compliance, and other high-stakes domains requiring verifiable provenance.

**Parametric Memory**   Compared with symbolic memory, parametric memory is *implicit, abstract*, and *generalizable*, making it naturally suited to tasks requiring conceptual understanding and broad pattern

induction. It is particularly effective when the agent must rely on general knowledge or rules that apply across diverse contexts, because such regularities can be internalized as distributed representations without requiring explicit external lookup. This internalization supports fluid reasoning and end-to-end processing, enabling the model to generalize systematically to unseen tasks or problem variations. Consequently, parametric memory is better aligned with tasks demanding structural insight, robust abstraction, and deeply ingrained behavioral or stylistic patterns.

**Possible Scenarios:**

- Role-playing or persona-consistent behaviors. (Shao et al., 2023; Zhou et al., 2024a)
- Mathematical reasoning, coding, games, and structured problem-solving.
- Human alignment and normative behavioral priors.
- Stylized, professional, or domain-expert responses.

**Latent Memory**  Unlike token-level or parametric memory, latent memory sits between explicit data and fixed model weights, enabling a unique balance of flexibility and efficiency. Its low readability provides intrinsic privacy protection, making latent representations suitable for sensitive information processing. At the same time, their high expressive capacity permits rich semantic encoding with minimal information loss, allowing agents to capture subtle correlations across modalities or tasks. Latent memory also supports efficient inference-time retrieval and integration, enabling agents to inject large quantities of compact knowledge. This memory type therefore prioritizes performance and scalability over interpretability, achieving high knowledge density and compression ideal for constrained or highly dynamic environments.
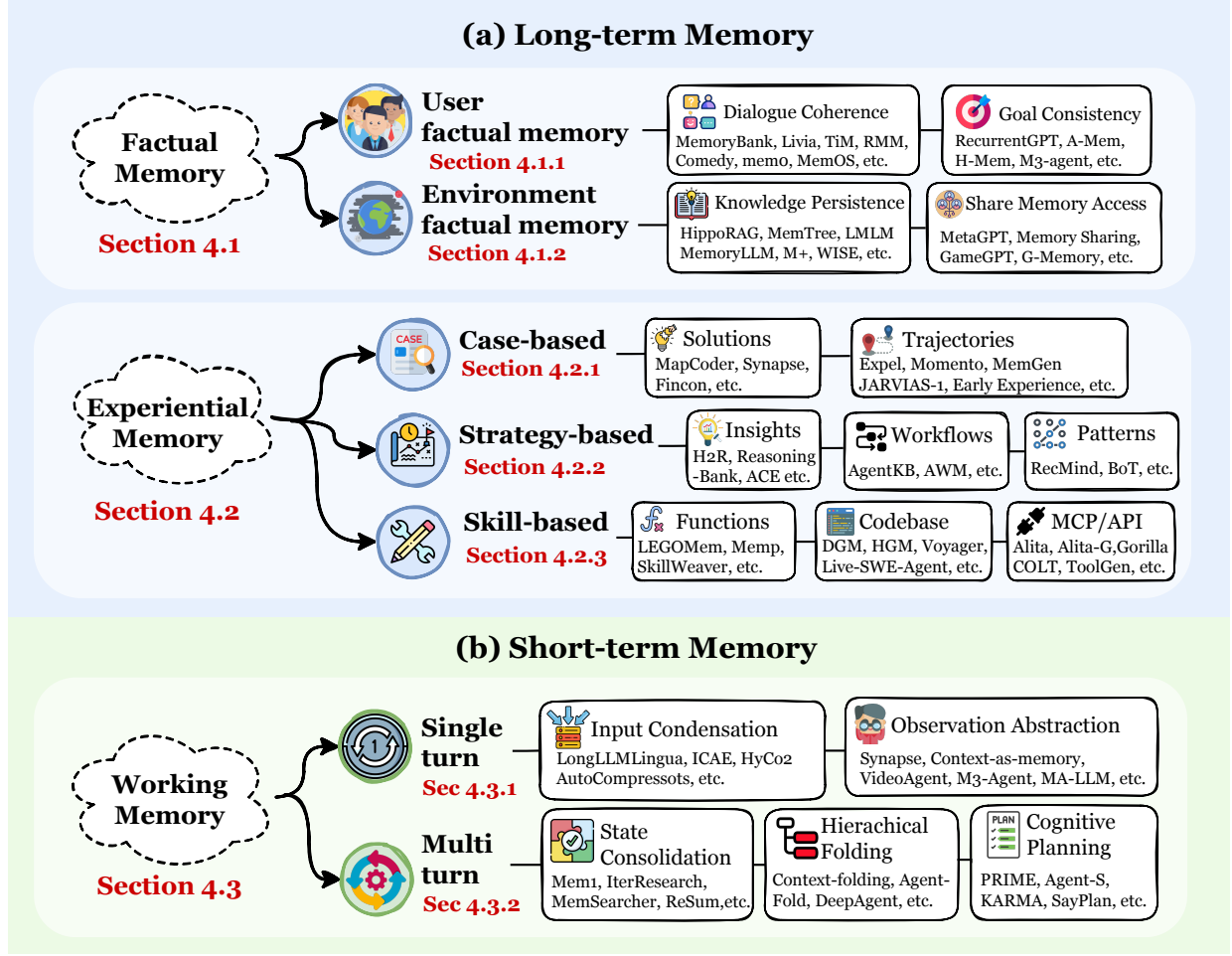
**Possible Scenarios:**

- Multimodal or fully integrated agent architectures. (Shi et al., 2025a; Cheng and Schwing, 2022; Wu et al., 2025d)
- On-device or edge deployment and cloud-serving environments.
- Encrypted or privacy-sensitive application domains.

## 4  Functions: Why Agents Need Memory?

The transition from large language models as general-purpose, stateless text processors to autonomous, goal-directed agents is not merely an incremental step but a fundamental paradigm shift. This shift exposes the critical limitation of statelessness. By definition, an agent must persist, adapt, and interact coherently over time. Achieving this relies not merely on a large context window but fundamentally on the capacity for **memory**. This section addresses the *functions*, or *fundamental purpose*, of agent memory, prioritizing the question of *why it is essential* over *how it is implemented*. We posit that agent memory is not a monolithic component but a set of distinct functional capabilities, each serving a unique objective in enabling persistent, intelligent behavior.

To provide a systematic analysis, this section organizes the *why* of memory around a functional taxonomy that maps directly to an agent's core requirements. At the highest level, we distinguish between two temporal categories: **long-term memory**, which serves as the persistent, cross-session store for accumulated knowledge, and **short-term memory**, which functions as the transient, in-session workspace for active reasoning. This high-level temporal split is further resolved into three primary functional pillars, which form the structure of our analysis. An overview of this taxonomy is provided in Figure 6.

**Figure 6** The functional taxonomy of agent memory. We organize memory capabilities based on their *functions* (purpose) into three primary pillars spanning two temporal domains: (1) **Factual Memory** serves as a persistent declarative knowledge base to ensure interaction *consistency*, *coherence*, and *adaptability*; (2) **Experiential Memory** encapsulates procedural knowledge to enable *continual learning* and *self-evolution* across episodes; and (3) **Working Memory** provides mechanisms for the active management of transient context.

---

**Three Primary Memory Functions**

1. **Factual Memory** (Section 4.1): The agent's declarative knowledge base, established to ensure consistency, coherence, and adaptability by recalling explicit facts, user preferences, and environmental states. This system answers the question: "What does the agent know?"

2. **Experiential Memory** (Section 4.2): The agent's procedural and strategic knowledge, accumulated to enable continual learning and self-evolution by abstracting from past trajectories, failures, and successes. This system answers: "How does the agent improve?"

3. **Working Memory** (Section 4.3): The agent's capacity-limited, dynamically controlled scratchpad for active context management during a single task or session. This system answers: "What is the agent thinking about now?"

---

These three memory systems are not isolated but form a dynamic, interconnected architecture that defines the agent's **cognitive loop**. The cycle begins with *encoding*, in which the outcomes of the agent's interactions, such as newly acquired facts or the results of a failed plan, are consolidated into long-term memory through

summarization, reflection, or abstraction. *Processing* subsequently occurs within working memory, which functions as the active workspace for immediate inference. To support this reasoning, the system relies on *retrieval* to populate the workspace with relevant context and skills drawn from the persistent stores of factual and experiential memory. This encoding-processing-retrieval sequence constitutes the central architectural pattern enabling agents to learn from the past simultaneously and reason in the present.

## 4.1 Factual Memory

Factual memory refers to the capacity of an agent to store and retrieve explicit, declarative **facts** about past events, user-specific information, and the state of the external environment. This information encompasses a wide range of content, including dialogue history, user preferences, and relevant properties of the external world. By allowing the agent to exploit historical information when interpreting current inputs, factual memory serves as the cornerstone for context awareness, personalized responses, and extended task planning.

To understand the structural composition of agent memory, we draw upon the cognitive science framework of *declarative memory* (Riedel and Blokland, 2015). In neuroscience, declarative memory denotes long-term storage for information that can be consciously accessed and is commonly analyzed in terms of two major components: *episodic* and *semantic* memory (Squire, 2004). *Episodic memory* stores personally experienced events associated with specific temporal and spatial contexts—the *what*, *where*, and *when* of an episode (Tulving, 1972, 2002). Its central characteristic is the capacity to mentally re-experience past events. *Semantic memory* retains general factual knowledge, concepts, and word meanings independent of the specific occasion on which they were acquired (Squire, 2004). While supported by a unitary declarative system in the human brain, these components represent distinct levels of abstraction.

In agent systems, this biological distinction is operationalized not as a rigid dichotomy but as a processing *continuum*. Systems typically initiate this process by logging concrete interaction histories as episodic traces, such as dialogue turns, user actions, and environment states (Zhong et al., 2024; Wang et al., 2024h; Chhikara et al., 2025). Subsequent processing stages apply summarization (Wang et al., 2025h; Chen et al., 2025c), reflection (Tan et al., 2025c; Park et al., 2023; Wang et al., 2025h), entity extraction (Gutierrez et al., 2024), and fact induction (Rasmussen et al., 2025). The resulting abstractions are stored in structures such as vector databases (Zhong et al., 2024), key-value stores, or knowledge graphs (Rasmussen et al., 2025; Sun et al., 2024), governed by procedures for deduplication and consistency checking. Through this sequence, raw event streams are gradually transformed into reusable semantic fact bases.

Functionally, this architecture ensures that the agent exhibits three fundamental properties during interaction: **consistency**, **coherence**, and **adaptability**.

- **Coherence** is reflected in robust context awareness. The agent can recall and integrate relevant interaction history, refer to past user inputs, and preserve topical continuity, ensuring responses form a logically connected dialogue rather than isolated utterances.

- **Consistency** implies stable behavior and self-presentation over time. By maintaining a persistent internal state regarding user-specific facts and its own commitments, the agent avoids contradictions and arbitrary changes of stance.

- **Adaptability** demonstrates the ability to personalize behavior based on stored user profiles and historical feedback. Consequently, response style and decision-making progressively align with the user's specific needs and characteristics.

For exposition, we further organize factual memory according to the primary entity it refers to. This entity-centric taxonomy, together with representative methods and their technical design choices, is systematically summarized in Table 4. This perspective highlights two central application domains:

> **Two Types of Factual Memory**
>
> - **User factual memory** (Section 4.1.1) denotes facts that sustain the consistency of interactions between humans and agents, including identities, stable preferences, task constraints, and historical commitments.
>
> - **Environment factual memory** (Section 4.1.2) denotes facts that sustain consistency with respect to the external world, such as document states, resource availability, and the capabilities of other agents.

### 4.1.1 User factual memory

User factual memory persists verifiable facts about a specific user across sessions and tasks, including identity, preferences, routines, historical commitments, and salient events.

Its primary function is to prevent characteristic failure modes of stateless interaction, such as coreference drift, repeated elicitation, and contradictory responses, thereby reducing interruptions to long-horizon goals (Tan et al., 2025c; Zhong et al., 2024). Engineering practice typically comprises selection and compression, structured organization, retrieval and reuse, and consistency governance, aiming to sustain *long-range dialogic and behavioral coherence* under bounded access cost.

**Dialogue Coherence**   Dialogue coherence requires an agent to preserve conversational context, user-specific facts, and a stable persona over extended periods. This ensures that later turns remain sensitive to earlier disclosures and affective cues, rather than degrading into repeated clarifications or inconsistent replies. To achieve this, modern systems implement user factual memory through two complementary strategies: *heuristic selection* and *semantic abstraction*.

To navigate finite context windows efficiently, a primary strategy is to *selectively* retain and rank interaction histories. Rather than retaining all raw logs, systems (Xi and Wang, 2025; Zhong et al., 2024; Park et al., 2023; Lei et al., 2025) maintain structured stores of past interactions, ranking entries by metrics such as *relevance*, *recency*, *importance*, or *distinctiveness*. By filtering retrieval based on these scores, high-value items are preserved and periodically condensed into higher-level summaries, conditioning subsequent responses to maintain continuity without overwhelming the agent's working memory.

Beyond mere selection, advanced frameworks emphasize the *transformation and abstraction* of raw dialogue fragments into higher-level semantic representations. Approaches such as Think in Memory (Liu et al., 2023a) and Reflective Memory Management (Tan et al., 2025c) convert raw interaction traces into *thought* representations or reflections via iterative update operations. This allows the agent to query a stable semantic memory, keeping later replies topically consistent and less repetitive. Similarly, COMEDY (Chen et al., 2025c) employs a single language model to generate, compress, and reuse memory while updating compact user profiles. These methods effectively stabilize *persona* and *preference* expression over long conversational histories by decoupling memory storage from the raw token surface form.

**Goal Consistency**   Goal consistency requires an agent to maintain and refine an explicit task representation over time. This ensures that clarifying questions, information requests, and actions remain strictly aligned with the primary objective, minimizing intent drift.

To mitigate such drift, systems utilize factual memory to dynamically track and update the task state. Approaches like RecurrentGPT (Zhou et al., 2023b), Memolet (Yen and Zhao, 2024), and MemGuide (Du et al., 2025b) retain confirmed information while highlighting unresolved elements. By guiding retrieval based on task intent, these methods help agents satisfy missing constraints and maintain focus across sessions.

For complex, long-horizon tasks, memory forms are often *structured* to facilitate localized retrieval centered on the active goal (Wu et al., 2025h). For instance, A-Mem (Xu et al., 2025c) organizes memories as an interconnected graph of linked notes, while H-Mem (Limbacher and Legenstein, 2020) employs associative mechanisms to recall prerequisite facts when subsequent steps depend on prior observations.

**Table 4** Taxonomy of factual memory methods. We categorize existing works based on the primary target entity: **User Factual Memory** focuses on sustaining interaction consistency, while **Environment Factual Memory** ensures consistency with the external world. Methods are compared across three technical dimensions: (1) **Carrier** (Section 3) identifies the storage medium, (2) **Structure** follows the taxonomy of token-level memory (Section 3.1), and (3) **Optimization** denotes the integration strategy, where *PE* encompasses prompt engineering and inference-time techniques without parameter updates, distinct from gradient-based methods like *SFT* and *RL*.

| Method | Carrier | Structure | Task | Optimization |
|---|---|---|---|---|
| *I. User factual Memory* | | | | |
| **(a) Dialogue Coherence** | | | | |
| MemGPT (Packer et al., 2023b) | Token-level | 1D | Long-term dialogue | PE |
| TiM (Park et al., 2023) | Token-level | 2D | QA | PE |
| MemoryBank (Zhong et al., 2024) | Token-level | 1D | Emotional Companion | PE |
| AI Persona (Wang et al., 2024f) | Token-level | 1D | Emotional Companion | PE |
| Encode-Store-Retrieve (Shen et al., 2024) | Token-level | 1D | Multimodal QA | PE |
| Livia (Xi and Wang, 2025) | Token-level | 1D | Emotional Companion | PE |
| mem0 (Chhikara et al., 2025) | Token-level | 1D | Long-term dialogue, QA | PE |
| RMM (Tan et al., 2025c) | Token-level | 2D | Personalization | PE, RL |
| D-SMART (Lei et al., 2025) | Token-level | 2D | Reasoning | PE |
| Comedy (Chen et al., 2025c) | Token-level | 1D | Summary, Compression, QA | PE |
| MEMENTO (Kwon et al., 2025) | Token-level | 1D | Embodied, Personalization | PE |
| O-Mem (Wang et al., 2025g) | Token-level | 3D | Personalized Dialogue | PE |
| DAM-LLM (Lu and Li, 2025) | Token-level | 1D | Emotional Companion | PE |
| MemInsight (Salama et al., 2025) | Token-level | 1D | Personalized Dialogue | PE |
| **(b) Goal Consistency** | | | | |
| RecurrentGPT (Zhou et al., 2023b) | Token-level | 1D | Long-Context Generation, Personalized Interactive Fiction | PE |
| Memolet (Yen and Zhao, 2024) | Token-level | 2D | QA, Document Reasoning | PE |
| MemGuide (Du et al., 2025b) | Token-level | 1D | Long-conv QA | PE, SFT |
| SGMem (Wu et al., 2025h) | Token-level | 2D | Long-context | PE |
| A-Mem (Xu et al., 2025c) | Token-level | 2D | QA, Reasoning | PE |
| M3-agent (Long et al., 2025) | Token-level | 2D | Multimodal QA | PE, SFT |
| *II. Environment factual Memory* | | | | |
| **(a) Knowledge Persistence** | | | | |
| MemGPT (Packer et al., 2023b) | Token-level | 1D | Document QA | PE |
| CALYPSO (Zhu et al., 2023) | Token-level | 1D | Tabletop Gaming | PE |
| AriGraph (Anokhin et al., 2024) | Token-level | 3D | Game, Multi-op QA | PE |
| HippoRAG (Gutierrez et al., 2024) | Token-level | 3D | QA | PE |
| WISE (Wang et al., 2024e) | Parametric | / | Document Reasoning, QA | SFT |
| MemoryLLM (Wang et al., 2024j) | Parametric | / | Document Reasoning | SFT |
| Zep (Rasmussen et al., 2025) | Token-level | 3D | Document analysis | PE |
| MemTree (Rezazadeh et al., 2025c) | Token-level | 2D | Document Reasoning, Dialogue | PE |
| LMLM (Zhao et al., 2025b) | Token-level | 1D | QA | SFT |
| M+ (Wang et al., 2025m) | Latent | / | Document Reasoning, QA | SFT |
| CAM (Li et al., 2025f) | Token-level | 3D | Multi-hop QA | SFT, RFT |
| MemAct (Zhang et al., 2025q) | Token-level | 1D | Multi-obj QA | RL |
| Mem-$\alpha$ (Wang et al., 2025o) | Token-Level | 1D | Document Reasoning | RL |
| WebWeaver (Li et al., 2025l) | Token-level | 1D | Deep Research | SFT |
| **(b) Shared Access** | | | | |
| GameGPT (Chen et al., 2023b) | Token-level | 1D | Game Development | PE |
| Generative Agent (Park et al., 2023) | Token-level | 2D | Social Simulation | PE |
| S³ (Gao et al., 2023a) | Token-level | 1D | Social Simulation | PE |
| Memory Sharing (Gao and Zhang, 2024a) | Token-level | 1D | Document Reasoning | PE |
| MetaGPT (Hong et al., 2024) | Token-level | 1D | Software Development | PE |
| G-Memory (Zhang et al., 2025e) | Token-level | 3D | QA | PE |
| OASIS (Yang et al., 2025) | Token-level, Parametric | 1D | Social Simulation | PE |

In embodied scenarios, factual memory grounds agent behavior in user-specific habits and environmental context. Systems such as M3-Agent (Long et al., 2025) and MEMENTO (Kwon et al., 2025) persist data on household members, object locations, and routines, reusing this information to minimize redundant exploration and repeated instructions. Similarly, Encode-Store-Retrieve (Shen et al., 2024) processes egocentric visual streams into text-addressable entries, allowing agents to answer questions based on past visual experiences without requiring user repetition.

**Summary**  Collectively, these mechanisms transform ephemeral interaction traces into a persistent cognitive substrate. By integrating retrieval-based ranking with generative abstraction, user factual memory upgrades the system from simple similarity matching to the active maintenance of explicit goals and constraints. This foundation yields a dual benefit: it fosters a sense of familiarity and trust through long-term behavioral coherence, while simultaneously enhancing operational efficiency by increasing task success rates, reducing redundancy, and lowering error recovery overhead.

### 4.1.2 Environment factual memory

Environment factual memory pertains to entities and states *external* to the user, encompassing long documents, codebases, tools, and interaction traces.

This memory paradigm addresses incomplete factual recall and unverifiable provenance, minimizes contradictions and redundancy in multi-agent collaboration, and stabilizes long-horizon tasks in heterogeneous environments. The central objective is to furnish an updatable, retrievable, and governable external fact layer, providing a stable reference across sessions and stages. Concretely, we categorize existing implementations along two complementary dimensions: *knowledge persistence* and *multi-agent shared access*.
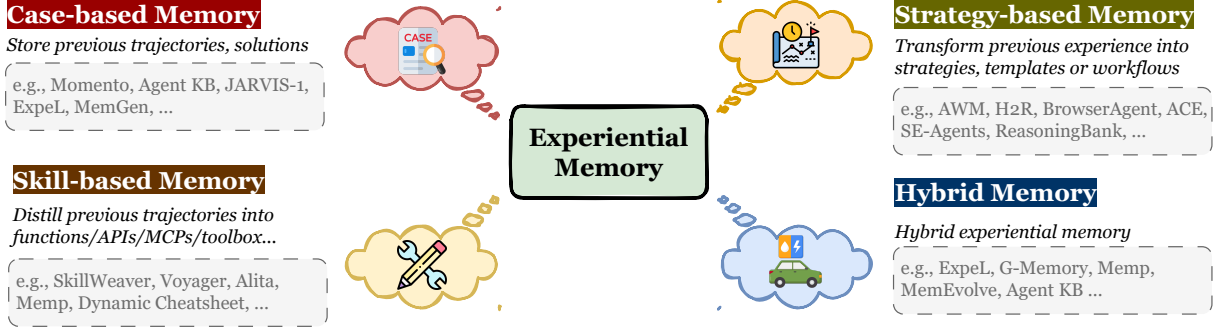
**Knowledge Persistence**  Knowledge memory refers to persistent representations of world knowledge and domain-specific knowledge that support long document analysis, factual question answering, multihop reasoning, and reliable retrieval of code and data resources.

In terms of *knowledge organization*, existing research focuses on structuring external data to enhance reasoning capabilities. For instance, HippoRAG (Gutierrez et al., 2024) utilizes knowledge graphs to facilitate evidence propagation, while MemTree (Rezazadeh et al., 2025c) employs a dynamic hierarchical structure to optimize aggregation and targeted access in growing corpora. Regarding storage form, LMLM (Zhao et al., 2025b) explicitly decouples factual knowledge from model weights by externalizing it into a database, thereby enabling direct knowledge edits and provenance verification without retraining. In narrative domains, CALYPSO (Zhu et al., 2023) distills lengthy game contexts into bite-sized prose, preserving critical story state accessibility.

In scenarios requiring continuous *knowledge updates*, parameter-centric approaches integrate persistence directly into the *model architecture*. Methods such as MEMORYLLM (Wang et al., 2024j), M+ (Wang et al., 2025m), and WISE (Wang et al., 2024e) incorporate trainable memory pools or side networks to absorb new information. Rather than relying solely on static external retrieval, these designs focus on the challenge of model editing, allowing agents to adapt to dynamic environments and correct obsolete facts while preserving the stability of the pre-trained backbone.

**Shared Access**  Shared memory establishes a visible and manageable common factual foundation for *multi-agent collaboration*, serving to align goals, carry intermediate artifacts, and eliminate redundant work. By maintaining a centralized repository of past queries and responses, frameworks such as Memory Sharing (Gao and Zhang, 2024b) enable agents to access and build on peers' accumulated insights asynchronously. This mechanism ensures that individual agents directly benefit from collective knowledge, thereby suppressing contradictory conclusions and enhancing overall system efficiency.

For complex project coordination, systems such as MetaGPT (Hong et al., 2024) and GameGPT (Chen et al., 2023b) utilize shared message pools as central workspaces for publishing plans and partial results. Similarly, G-Memory (Zhang et al., 2025e) employs hierarchical memory graphs as a unified coordination medium. These architectures facilitate consistency maintenance around the current project state, which reduces communication overhead and enables the extraction of reusable workflows from historical collaborations.

**Figure 7** Taxonomy of experiential memory paradigms. We classify approaches based on the *abstraction level* of stored knowledge: (1) **Case-based Memory** preserves raw trajectories and solutions as concrete exemplars; (2) **Strategy-based Memory** abstracts experiences into high-level strategies, templates, or workflows; (3) **Skill-based Memory** distills procedural knowledge into executable functions and APIs; and (4) **Hybrid Memory** integrates multiple representations. Together, these systems mirror human *procedural memory* to enable continual learning and self-evolution. This figure draws inspiration from Gao et al. (2025).

In the domain of social simulation, platforms like Generative Agents (Park et al., 2023) and S$^3$ (Gao et al., 2023a), alongside large-scale simulators such as OASIS (Yang et al., 2025) and AgentSociety (Piao et al., 2025), model the global environment and public interaction logs as a shared memory substrate. This substrate is incrementally updated and observed by the population, allowing information to diffuse naturally among agents and supporting coherent, history-aware social dynamics at scale.

**Summary** environment factual memory furnishes a continuously updatable, auditable, and reusable external fact layer. On the knowledge axis, it improves completeness, interpretability, and editability of factual recall through structured organization and long-term memory modules. On the collaboration axis, it maintains cross-agent and cross-stage consistency through sharing and governance, thereby enabling robust decision-making and execution under long horizons, multiple actors, and multi-source information.

## 4.2 Experiential Memory

Experiential memory encapsulates the mechanism by which agents encode historical trajectories, distilled strategies, and interaction outcomes into durable, retrievable representations. Unlike working memory, which manages transient context, experiential memory focuses on the long-term accumulation and transfer of knowledge *across* distinct episodes.

Theoretically grounded in cognitive science, this paradigm parallels human *nondeclarative memory*, specifically the *procedural* and *habit* systems (Squire, 2004; Seger and Spiering, 2011). Biological systems rely on distributed neural circuits for implicit skill acquisition (Reber, 2013). In contrast, agentic experiential memory typically employs explicit data structures, such as vector databases or symbolic logs. This implementation difference grants agents a unique capability absent in biological counterparts: the ability to introspect, edit, and reason over their own procedural knowledge.

Crucially, experiential memory serves as a foundation for **continual learning** and **self-evolution** in the *era of experience* (Sutton, 2025; Gao et al., 2025). By maintaining a repository of structured experiences, agents achieve a non-parametric path to adaptation and avoid the prohibitive costs of frequent parametric updates. This mechanism effectively closes the learning loop by converting interaction feedback into reusable knowledge. Through this process, agents rectify past errors, abstract generalizable heuristics, and compile routine behaviors. Consequently, such adaptation minimizes redundant computations and refines decision-making over time (Zhao et al., 2024; Shinn et al., 2023b).

To systematically analyze existing literature, we classify experiential memory based on the *abstraction level* of the stored information. An overview of this abstraction-based taxonomy and representative paradigms is illustrated in Figure 7. Representative methods under this abstraction-based taxonomy, together with their storage carriers, representation forms, and optimization strategies, are summarized in Table 5.

> **Three Types of Experiential Memory**
>
> - **Case-based Memory** (Section 4.2.1) stores minimally processed records of historical episodes, prioritizing high informational fidelity to support direct replay and imitation. By retaining the original alignment between situations and outcomes, it serves as a repository of concrete, verifiable evidence that functions as in-context exemplars for evidence-driven learning.
>
> - **Strategy-based Memory** (Section 4.2.2) distills transferable reasoning patterns, workflows, and high-level insights from past trajectories to guide planning across diverse scenarios. Acting as a cognitive scaffold, it decouples decision-making logic from specific contexts, thereby enhancing cross-task generalization and constraining the search space for complex reasoning.
>
> - **Skill-based Memory** (Section 4.2.3) encapsulates executable procedural capacities, ranging from atomic code snippets to standardized API protocols, that operationalize abstract strategies into verifiable actions. This category serves as the agent's active execution substrate, enabling the modular expansion of capabilities and the efficient handling of tool-use environments.

**Table 5** Taxonomy of experiential memory methods. We categorize existing works based on the *abstraction level* of stored knowledge: **Case-based Memory** preserves raw records for direct replay, **Strategy-based Memory** distills abstract heuristics for planning, and **Skill-based Memory** compiles executable capabilities for action. Methods are compared across three technical dimensions: (1) **Carrier** (Section 3) identifies the storage medium, (2) **Form** specifies the representation format of the experience, and (3) **Optimization** denotes the integration strategy, where *PE* encompasses prompt engineering and inference-time techniques without parameter updates, distinct from gradient-based methods like *SFT* and *RL*.

| Method | Carrier | Form | Task | Optimization |
|---|---|---|---|---|
| *I. Case-based Memory* | | | | |
| Expel (Zhao et al., 2024) | Token-level | Solution | Reasoning | PE |
| Synapse (Zheng et al., 2024a) | Token-level | Solution | Web Interaction, Instruction-guided Web Task | PE |
| Fincon (Yu et al., 2024) | Token-level | Solution | Financial | PE |
| MapCoder (Islam et al., 2024) | Token-level | Solution | Coding | PE |
| Memento (Zhou et al., 2025a) | Token-level | Trajectory | Reasoning | RL |
| COLA (Zhao et al., 2025a) | Token-level | Trajectory | GUI, Web Navigation, Reasoning | PE |
| Continuous Memory (Wu et al., 2025e) | Latent | Trajectory | GUI | SFT |
| JARVIS-1 (Wang et al., 2025p) | Token-level | Trajectory | Game, GUI Interaction | PE |
| MemGen (Zhang et al., 2025d) | Latent | Trajectory | Web Search, Embodied Simulation, Reasoning, Math, Code | RL, SFT |
| Early Experience (Zhang et al., 2025j) | Parametric | Trajectory | Embodied Simulation, Reasoning, Web Navigation | SFT |
| DreamGym (Chen et al., 2025e) | Token-level | Trajectory | Web Interaction, Embodied Simulation, Shopping | RL |
| *II. Strategy-based Memory* | | | | |
| Reflexion (Shinn et al., 2023a) | Token-level | Insight | Embodied Simulation, Reasoning, Coding | PE |
| Buffer of Thoughts (Yang et al., 2024b) | Token-level | Pattern | Game, Reasoning, Coding | PE |
| AWM (Wang et al., 2024l) | Token-level | Workflow | Web Interaction, Instruction-guided Web Task | PE |
| RecMind (Wang et al., 2024h) | Token-level | Pattern | Recommendation | PE |
| H$^2$R (Ye et al., 2025b) | Token-level | Insight | Game, Embodied Simulation | PE |
| ReasoningBank (Ouyang et al., 2025) | Token-level | Insight | Web Interaction, Instruction-guided Web Task | PE |
| R2D2 (Huang et al., 2025c) | Token-level | Insight | Web Interaction | PE |
| BrowserAgent (Yu et al., 2025d) | Token-level | Insight | General QA, Web search | RL, SFT |
| Agent KB (Tang et al., 2025d) | Token-level | Workflow | Code, Reasoning | PE |

<div align="right">Continued on next page</div>

**Table 5** Taxonomy of experiential memory methods. We categorize existing works based on the *abstraction level* of stored knowledge: **Case-based Memory** preserves raw records for direct replay, **Strategy-based Memory** distills abstract heuristics for planning, and **Skill-based Memory** compiles executable capabilities for action. (continued)

| Method | Carrier | Form | Task | Optimization |
|---|---|---|---|---|
| ToolMem (Xiao et al., 2025b) | Token-level | Insight | Reasoning, Image Generation | PE |
| PRINCIPLES (Kim et al., 2025a) | Token-level | Pattern | Emotional Companion | PE |
| SE-Agent (Sun et al., 2025b) | Token-level | Insight | Coding | PE |
| ACE (Zhang et al., 2025m) | Token-level | Insight | Coding, Tool calling, Financial | PE |
| Flex (Cai et al., 2025b) | Token-level | Insight | Math, Chemistry, Biology | PE |
| AgentEvolver (Zhai et al., 2025) | Parametric | Pattern | Tool-augmented Task | RL |
| Dynamic Cheatsheet (Suzgun et al., 2025) | Token-level | Insight | Math, Reasoning, Game | PE |
| Training-Free GRPO (Cai et al., 2025a) | Token-level | Insight | Math, Reasoning, Web Search | PE |
| *III. Skill-based Memory* | | | | |
| CREATOR (Qian et al., 2023) | Token-level | Function and Script | Reasoning, Math | PE |
| Gorilla (Patil et al., 2024) | Token-level | API | Tool calling | SFT |
| ToolRerank (Zheng et al., 2024b) | Token-level | API | Tool calling | PE |
| Voyager (Wang et al., 2024b) | Token-level | Code Snippet | Game | PE |
| RepairAgent (Bouzenia et al., 2024) | Token-level | Function and Script | Coding | PE |
| COLT (Qu et al., 2024) | Token-level | API | Tool calling | SFT |
| ToolLLM (Qin et al., 2024a) | Token-level | API | Tool Calling | SFT |
| LEGOMem (Han et al., 2025a) | Token-level | Function and Script | Office | PE |
| Darwin Gödel Machine (Zhang et al., 2025h) | Token-level | Code Snippet | Code | PE |
| Huxley-Gödel Machine (Wang et al., 2025j) | Token-level | Code Snippet | Code | PE |
| Memp$^P$ (Fang et al., 2025d) | Token-level | Function and Script | Embodied Simulation, Travel Planning | PE |
| SkillWeaver (Zheng et al., 2025a) | Token-level | Function and Script | Web Interaction, Instruction-guided Web Task | PE |
| Alita (Qiu et al., 2025c) | Token-level | MCP | Math, Reasoning, VQA | PE |
| Alita-G (Qiu et al., 2025b) | Token-level | MCP | Math, Reasoning, VQA | PE |
| LearnAct (Liu et al., 2025a) | Token-level | Function and Script | Mobile GUI | PE |
| ToolGen (Wang et al., 2025i) | Parametric | API | Tool calling | SFT |
| MemTool (Lumer et al., 2025) | Token-level | MCP | Tool calling | SFT |
| ToolRet (Shi et al., 2025c) | Token-level | API | Web, Code, Tool Retrieval | SFT |
| DRAFT (Qu et al., 2025a) | Token-level | API | Tool calling | PE |
| ASI (Wang et al., 2025r) | Token-level | Functions and Scripts | Web Interaction | PE |

### 4.2.1 Case-based Memory

Case-based memory stores minimally processed records of historical *events*, prioritizing fidelity to ensure that episodes can be replayed or reused as in-context exemplars. Unlike strategy templates or skill modules, cases avoid extensive abstraction, thereby preserving the original alignment between situations and solutions.

**Trajectories** This category preserves interaction sequences to enable replay and evidence-driven learning. To optimize retrieval in text-based environments, Memento (Zhou et al., 2025a) employs soft Q-learning to dynamically refine the probability of selecting high-utility past trajectories. In multimodal settings, JARVIS-1 (Wang et al., 2025p), EvoVLA (Liu et al., 2025h) and Auto-scaling Continuous Memory (Wu et al., 2025e) retain visual context, with the former storing survival experiences in Minecraft and the latter compressing GUI history into continuous embeddings. Furthermore, the early experience paradigm (Zhang et al., 2025j) constructs reward-free, agent-generated interaction traces and integrates them into model parameters via mid-training to enhance generalization.

**Solutions** This category treats memory as a repository of proven solutions. ExpeL (Zhao et al., 2024) autonomously gathers experience through trial-and-error, storing successful trajectories as exemplars while

extracting textual insights to guide future actions. Synapse (Zheng et al., 2024a) similarly injects abstracted state-action episodes as contextual examples to align problem-solving patterns. In program synthesis, MapCoder (Islam et al., 2024) keeps relevant example code as a playbook-like case that multi-agent pipelines retrieve and adapt to improve reliability on complex tasks. In the financial domain, FinCon (Yu et al., 2024) maintains an episodic memory of past actions, PnL trajectories, and belief updates to facilitate robust cross-round decision-making.

**Summary**  Case-based memory offers high informational fidelity and provides verifiable evidence for imitation. However, the reliance on raw data imposes challenges regarding retrieval efficiency and context window consumption. Distinguished from executable skills or abstract strategies, cases do not encompass orchestration logic or function interfaces. Instead, they serve as the factual substrate upon which higher-level reasoning operates.

### 4.2.2  Strategy-based Memory

Unlike case libraries that retain *what happened*, strategy-based memory extracts transferable knowledge of *how to act*, encompassing reusable reasoning patterns, task decompositions, insights, abstractions, and cross-situational workflows. It elevates experiences into editable, auditable, and composable high-level knowledge, thereby reducing dependence on lengthy trajectory replay and improving cross-task generalization and efficiency. We focus on non-code or weakly code-based templates and workflows in this section, while executable functions, APIs, MCP protocols, and code snippets are classified under Section 4.2.3. Based on the granularity and structural complexity of the retained knowledge, we categorize strategy-based memory into three distinct types: atomic **Insights**, sequential **Workflows**, and schematic **Patterns**.

**Insights**  This category of approaches focuses on distilling discrete pieces of knowledge, such as granular decision *rules* and reflective *heuristics*, from past trajectories. $H^2R$ (Ye et al., 2025b) explicitly decouples planning-level and execution-level memories, enabling high-level planning insights and low-level operational rules to be retrieved separately for fine-grained transfer in multi-task scenarios. R2D2 (Huang et al., 2025c) integrates remembering, reflecting, and dynamic decision-making for web navigation, deriving corrective insights from both failed and successful cases to inform subsequent episodes. For long-horizon web automation, BrowserAgent (Yu et al., 2025d) persists key conclusions as explicit memory to stabilize extended chains of reasoning and mitigate context drift.

**Workflows**  Distinct from atomic, static insights, *workflows* encapsulate strategies as structured *sequences* of actions—executable routines abstracted from prior trajectories to guide multi-step execution at inference time. Agent Workflow Memory (AWM) (Wang et al., 2024l) induces reusable workflows on Mind2Web (Deng et al., 2023) and WebArena (Zhou et al., 2023a) and uses them as high-level scaffolds to guide subsequent generation, improving success rates and reducing steps without updating base model weights. This demonstrates that strategy templates can act as a top-level controller that complements case-level evidence. Agent KB (Tang et al., 2025d) establishes a unified knowledge base that treats workflows as transferable procedural knowledge. It employs hierarchical retrieval, accessing workflows first to structure the strategic approach and enabling problem-solving logic reuse across diverse agent architectures.

**Patterns**  At a higher level of abstraction, reasoning patterns function as *cognitive templates* that encapsulate the structure of problem-solving, enabling agents to tackle complex reasoning tasks by instantiating these *generalizable skeletons*. Buffer of Thoughts (Yang et al., 2024b) maintains a meta-buffer of thought templates that are retrieved and instantiated to solve new problems. Similarly, ReasoningBank (Ouyang et al., 2025) abstracts both successes and failures into reusable reasoning units, facilitating test-time expansion and robust learning. RecMind's self-inspiring planning algorithm (Wang et al., 2024h) generates intermediate self-guidance to structure subsequent planning and tool use. In the domain of dialogue agents, PRINCIPLES (Kim et al., 2025a) builds a synthetic strategy memory via offline self-play to guide strategy planning at inference, thereby eliminating the need for additional training. These advances indicate a paradigmatic shift from descriptive rules to portable reasoning structures.

**Summary**   Strategy-based memory, which encompasses insights, workflows, and patterns, serves as a high-level scaffold to guide generative reasoning. Unlike case-based memory that relies on retrieving specific, raw trajectories which may be noisy or context-dependent, this form of memory distills generalizable schemas to effectively constrain the search space and improve robustness on unseen tasks. However, a key distinction is that these strategies function as structural guidelines rather than executable actions; they direct the planning process but do not interact with the environment directly. This limitation necessitates skill-based memory, discussed in the following section, which stores callable capabilities and tools. Ultimately, robust agents typically synergize these components: strategies provide the abstract planning logic, while skills handle the grounded execution.

### 4.2.3 Skill-based Memory

Skill memory captures an agent's procedural capacity and operationalizes abstract strategy into verifiable actions. It encodes what the agent can do, complements declarative knowledge of what the agent knows, and anchors the perception–reasoning–action loop by providing invocable, testable, and composable executables. Recent evidence shows that language models can learn when and how to call tools and scale reliably with large tool repertoires, establishing skill memory as the execution substrate of modern agents.

Skill memory spans a continuum from internal, fine-grained code to externalized, standardized interfaces. The unifying criteria are straightforward: skills must be **callable** by the agent, their outcomes must be verifiable to support learning, and they must compose with other skills to form larger routines.

**Code Snippets**   Executable code stored as reusable snippets offers the fastest path from experience to capability. In open-ended tasks, agents distill successful sub-trajectories into interpretable programs and reuse them across environments. Voyager (Wang et al., 2024b) exemplifies this pattern with an ever-growing skill library; the Darwin Gödel Machine (Zhang et al., 2025h) goes further by safely rewriting its own code under empirical validation, yielding self-referential and progressively more capable skill sets.

**Functions and Scripts**   Abstracting complex behaviors into modular functions or scripts enhances reusability and generalization. Recent advancements empower agents to autonomously create specialized tools for problem-solving (Qian et al., 2023; Yuan et al., 2024a), and to refine tool-use capabilities through demonstrations and environmental feedback across diverse domains such as mobile GUIs, web navigation, and software engineering (Fang et al., 2025d; Zheng et al., 2025a; Bouzenia et al., 2024). Furthermore, emergent mechanisms for procedural memory enable agents to distill execution trajectories into retrievable scripts, facilitating efficient generalization to novel scenarios (Liu et al., 2025a; Han et al., 2025a).

**APIs**   APIs serve as the universal interface for encapsulated skills. While earlier work focused on fine-tuning models to correctly invoke tools (Schick et al., 2023; Patil et al., 2024), the exponential growth of API libraries has shifted the primary bottleneck to retrieval. Standard information retrieval methods often fail to capture the functional semantics of tools (Shi et al., 2025c). Consequently, recent approaches have moved towards learning-based retrieval and reranking strategies that account for tool documentation quality, hierarchical relationships, and collaborative usage patterns to bridge the gap between user intent and executable functions (Zheng et al., 2024b; Gao and Zhang, 2024c; Qu et al., 2024, 2025a).

**MCPs**   To reduce protocol fragmentation in API-based ecosystems, the Model Context Protocol provides an open standard that unifies how agents discover and use tools and data, including code-execution patterns that load tools on demand and cut context overhead (Qiu et al., 2025c,b). Broad platform support indicates a convergence toward a common interface layer.

Beyond standard executables, research explores learnable memories of tool capabilities to handle uncertain neural tools, parametric integration that embeds tool symbols to unify retrieval and calling, and architecture-as-skill perspectives where specialized agents are callable modules within a modular design space (Xiao et al., 2025b; Wang et al., 2025i; Zhao et al., 2025a). Collectively, these strands reframe skill memory as a learnable, evolving, and orchestrable capability layer.

**Summary** In conclusion, skill-based memory constitutes the active execution substrate of the agent, evolving from static code snippets and modular scripts to standardized APIs and learnable architectures. It bridges the gap between abstract planning and environmental interaction by operationalizing insights from case-based and strategy-based memories into verifiable procedures. As mechanisms for tool creation, retrieval, and interoperability (e.g., MCP) mature, skill memory moves beyond simple storage, enabling a continuous loop of capability synthesis, refinement, and execution that drives open-ended agent evolution.

### 4.2.4 Hybrid memory

Advanced agent architectures increasingly adopt a **hybrid** design that integrates multiple forms of experiential memory to balance grounded evidence with generalizable logic. By maintaining a spectrum of knowledge spanning raw episodes, distilled rules, and executable skills, these systems dynamically select the most appropriate memory format, ensuring both retrieval precision and broad generalization across contexts.

A prominent direction involves coupling *case-based* and *strategy-based* memories to facilitate complementary reasoning. For example, ExpeL (Zhao et al., 2024) synergizes concrete trajectories with abstract textual insights, allowing agents to recall specific solutions while applying general heuristics. Agent KB (Tang et al., 2025d) employs a hierarchical structure where high-level workflows guide planning and specific solution paths provide execution details. Similarly, R2D2 (Huang et al., 2025c) integrates a replay buffer of historical traces with a reflective mechanism that refines decision strategies from past errors, effectively bridging case retrieval and strategic abstraction. Complementing these, Dynamic Cheatsheet (Suzgun et al., 2025) prevents redundant computation by storing accumulated strategies and problem-solving insights for immediate reuse at inference time.

Furthermore, recent frameworks strive to unify the lifecycle of memory, incorporating Skill-based components or establishing comprehensive cognitive architectures. In scientific reasoning, ChemAgent (Tang et al., 2025c) constructs a self-updating library that pairs execution cases with decomposable skill modules, enabling the model to refine its chemical reasoning through accumulated experience. Taking a holistic approach, LARP (Yan et al., 2023) establishes a cognitive architecture for open-world games that harmonizes semantic memory for world knowledge, episodic memory for interaction cases, and procedural memory for learnable skills, ensuring consistent role-playing and robust decision-making. Finally, evolutionary systems like G-Memory (Zhang et al., 2025c) and Memp (Fang et al., 2025d) implement dynamic transitions, where repeated successful cases are gradually compiled into efficient skills, automating the shift from heavy retrieval to rapid execution. A recent effort, MemVerse (Liu et al., 2025d) combines both parametric memory and token-level prcedural memory.

## 4.3 Working Memory

In cognitive science, working memory is defined as a capacity-limited, dynamically controlled mechanism that supports higher-order cognition by selecting, maintaining, and transforming task-relevant information in the moment (Baddeley, 2012). Beyond mere temporary storage, it implies active control under resource constraints. This perspective is grounded in frameworks such as the multicomponent model and the embedded-processes account, both of which emphasize attentional focus, interference control, and bounded capacity (Cowan, 2014).

When transposed to LLMs, the standard context window functions primarily as a passive, read-only buffer. Although the model can consume the window's contents during inference, it lacks explicit mechanisms to select, sustain, or transform the current workspace dynamically. Recent behavioral evidence suggests that current models do not exhibit human-like working memory characteristics, underscoring the necessity for explicitly engineered, operable working memory mechanisms (Huang et al., 2025a).

Throughout this section, we define working memory as the set of mechanisms for the **active management and manipulation** of context within a single episode (Zhang et al., 2025q). The objective is to transform the context window from a passive buffer into a controllable, updatable, and interference-resistant workspace. This transition offers immediate benefits: it increases the density of task-relevant information under fixed attention budgets, suppresses redundancy and noise, and enables the rewriting or compression of representations to preserve coherent chains of thought. We categorize these mechanisms based on the interaction dynamics. Representative working memory approaches under this interaction-based taxonomy, together with their storage carriers, task domains, and optimization strategies, are systematically summarized in Table 6.

> **Two Types of Working Memory**
>
> - **Single-turn Working Memory** (Section 4.3.1) focuses on **input condensation and abstraction**. In this setting, the system must process massive immediate inputs such as long documents or high-dimensional multimodal streams within a single forward pass. The goal is to dynamically filter and rewrite evidence to construct a bounded computational scratchpad, thereby maximizing the effective information payload per token.
>
> - **Multi-turn Working Memory** (Section 4.3.2) addresses **temporal state maintenance**. In sequential interactions, the challenge is to prevent historical accumulation from overwhelming the attention mechanism. This involves maintaining task states, goals, and constraints through a continuous loop of reading, executing, and updating, ensuring that intermediate artifacts are folded and consolidated across turns.

In summary, working memory for LLMs represents a paradigm shift towards active, within-episode context management. By aligning with the cognitive requirement of active manipulation, it suppresses interference and provides a practical solution to the engineering constraints of long-context inference.

### 4.3.1 Single-turn Working Memory

Single-turn working memory addresses the challenge of processing massive immediate inputs, including long documents (Chevalier et al., 2023) and high-dimensional multimodal streams (Wang et al., 2024g), within a single forward pass. Rather than passively consuming the entire context, the objective is to actively construct a writable workspace. This involves filtering and transforming raw information to increase density and operability under fixed attention and memory budgets (Jiang et al., 2023, 2024). We categorize these mechanisms into *input condensation*, which reduces physical token count, and *observation abstraction*, which transforms data into structured semantic representations.

**Input Condensation** Input condensation techniques aim to preprocess the context to minimize token usage while preserving essential information (Jiang et al., 2023). These methods generally fall into three paradigms: hard, soft, and hybrid condensation (Liao et al., 2025a).

*Hard condensation* discretely selects tokens based on importance metrics. Approaches like LLMLingua (Jiang et al., 2023) and LongLLMLingua (Jiang et al., 2024) estimate token perplexity to discard predictable or task-irrelevant content, while CompAct (Yoon et al., 2024) adopts an iterative strategy to retain segments that maximize information gain. Although efficient, hard selection risks severing syntactic or semantic dependencies. *Soft condensation* encodes variable-length contexts into dense latent vectors (memory slots). Methods such as Gist (Mu et al., 2023), In-Context Autoencoder (ICAE) (Ge et al., 2024), and AutoCompressors (Chevalier et al., 2023) train models to compress prompts into valid summary tokens or distinct memory embeddings. This achieves high compression ratios but requires additional training and may obscure fine-grained details. *Hybrid* approaches like HyCo2 (Liao et al., 2025a) attempt to reconcile these *trade-offs* by combining global semantic adapters (soft) with token-level retention probabilities (hard).

**Observation Abstraction** While condensation focuses on reduction, *observation abstraction* aims to transform raw observations into structured formats that facilitate reasoning. This mechanism maps dynamic, high-dimensional observation spaces into fixed-size memory states, preventing agents from being overwhelmed by raw data.

In complex interactive environments, abstraction converts verbose inputs into concise state descriptions. Synapse (Zheng et al., 2024a) rewrites unstructured HTML DOM trees into task-relevant state summaries to guide GUI automation. Similarly, in multimodal settings, processing every frame of a video stream is computationally prohibitive. Working memory mechanisms address this by extracting semantic structures: Context as Memory (Yu et al., 2025b) filters frames based on field-of-view overlap, VideoAgent (Wang et al., 2024g) converts streams into temporal event descriptions, and MA-LMM (He et al., 2024) maintains a bank of visual features. These methods effectively rewrite high-dimensional, redundant streams into low-dimensional, semantically rich representations operable within a limited context window for efficient processing.

**Table 6** Taxonomy of working memory methods. We categorize approaches into **Single-turn** and **Multi-turn** settings based on interaction dynamics. Methods are compared across three technical dimensions: (1) **Carrier** (Section 3) identifies the storage medium, (2) **Task** specifies the evaluation domain or application scenario, and (3) **Optimization** denotes the integration strategy, where PE encompasses prompt engineering and inference-time techniques without parameter updates, distinct from gradient-based methods like SFT and RL.

| Method | Carrier | Task | Optimization |
|---|---|---|---|
| *I. Single-turn Working Memory* | | | |
| **(a) Input Condensation** | | | |
| Gist (Mu et al., 2023) | Latent | Instruction Fine-tuning | SFT |
| ICAE (Ge et al., 2024) | Latent | Language Modeling, Instruction Fine-tuning | Pretrain, LoRA |
| AutoCompressors (Chevalier et al., 2023) | Latent | Langague Modeling | SFT |
| LLMLingua (Jiang et al., 2023) | Token-level | Reasoning, Conversation, Summarization | PE |
| LongLLMLingua (Jiang et al., 2024) | Token-level | Multi-doc QA, Long-context, Multi-hop QA | PE |
| CompAct (Yoon et al., 2024) | Token-level | Document QA | SFT |
| HyCo2 (Liao et al., 2025a) | Hybrid | Summarization, Open-domain QA, Multi-hop QA | SFT |
| Sentence-Anchor (Tarasov et al., 2025) | Latent | Document QA | SFT |
| MELODI (Chen et al., 2024c) | Hybrid | Pretraining | Pretrain |
| **(b) Observation Abstraction** | | | |
| Synapse (Zheng et al., 2024a) | Token-level | Computer Control, Web Navigation | PE |
| VideoAgent (Wang et al., 2024g) | Token-level | Long-term Video Understanding | PE |
| MA-LMM (He et al., 2024) | Latent | Long-term Video Understanding | SFT |
| Context as Memory (Yu et al., 2025b) | Token-level | Long-term Video Generation | PE |
| *II. Multi-turn Working Memory* | | | |
| **(c) State Consolidation** | | | |
| MEM1 (Zhou et al., 2025b) | Latent | Retrieval, Open-domain QA, Shopping | RL |
| MemGen (Zhang et al., 2025d) | Latent | Reasoning, Embodied Action, Web Search, Coding | RL |
| MemAgent (Yu et al., 2025a) | Token-level | Long-term Doc. QA | RL |
| ReMemAgent (Shi et al., 2025b) | Token-level | Long-term Doc. QA | RL |
| ReSum (Wu et al., 2025f) | Token-level | Long-horizon Web Search | RL |
| MemSearcher (Yuan et al., 2025a) | Token-level | Multi-hop QA | SFT, RL |
| ACON (Kang et al., 2025c) | Token-level | App use, Multi-objective QA | PE |
| IterResearch (Chen et al., 2025a) | Token-level | Reasoning, Web Navigation, Long-Horizon QA | RL |
| SUPO (Lu et al., 2025a) | Token-level | Long-horizon task | RL |
| AgentDiet (Xiao et al., 2025a) | Token-level | Long-horizon task | PE |
| SUMER (Zheng et al., 2025c) | Token-level | QA | RL |
| **(d) Hierarchical Folding** | | | |
| HiAgent (Hu et al., 2025a) | Token-level | Long-horizon Agent Task | PE |
| Context-Folding (Zhang et al., 2025q) | Token-level | Deep Research, SWE | RL |
| AgentFold (Ye et al., 2025a) | Token-level | Web Search | SFT |
| DeepAgent (Li et al., 2025h) | Token-level | Tool Use, Shopping, Reasoning | RL |
| **(e) Cognitive Planning** | | | |
| SayPlan (Rana et al., 2023) | Token-level | 3D Scene Graph, Robotics | PE |
| KARMA (Wang et al., 2025q) | Token-level | Household | PE |
| Agent-S (Agashe et al., 2025) | Token-level | Computer Use | PE |
| PRIME (Tran et al., 2025) | Token-level | Multi-hop QA, Knowledge-intensive Reasoning | PE |

**Summary**  Single-turn working memory functions as an *active compression layer* that maximizes the utility of the context window for immediate reasoning. By employing input condensation and observation abstraction, these mechanisms effectively increase the information density of the operational workspace, ensuring that critical evidence is retained despite capacity constraints. However, this optimization is strictly *intra-turn*; it addresses the breadth and complexity of static inputs rather than the temporal continuity of dynamic interactions.

### 4.3.2 Multi-turn Working Memory

Multi-turn working memory addresses a fundamentally different problem space than the single-turn setting. In long-horizon interactions, the primary bottleneck shifts from instantaneous context capacity to the continuous maintenance of **task state** and **historical relevance**. Even with extended context windows, the accumulation of history inevitably saturates attention budgets, increases latency, and induces goal drift (Lu et al., 2025b). To mitigate this, working memory in multi-turn settings functions as an externalized state carrier, organizing a continuous loop of reading, evaluation, and writing. The objective is to preserve critical state information accessible and consistent within a bounded resource budget. We categorize these mechanisms by their state management strategies: *state consolidation*, *hierarchical folding*, and *cognitive planning*.

**State Consolidation**    In continuous interaction streams, state consolidation maps an ever-growing trajectory into a fixed-size state space through dynamic updates. Treating interaction as a streaming environment, MemAgent (Yu et al., 2025a), and MemSearcher (Yuan et al., 2025a) employ recurrent mechanisms to update fixed-budget memory and discard redundancy, answering queries from a compact, evolving state. ReSum (Wu et al., 2025f) extends this by periodically distilling history into reasoning states, utilizing reinforcement learning to optimize summary-conditioned behavior for indefinite exploration.

Beyond heuristic summarization, ACON (Kang et al., 2025c) frames state consolidation as an optimization problem, jointly compressing environment observations and interaction histories into a bounded condensation and iteratively refining compression guidelines from failure cases. IterResearch (Chen et al., 2025a) further adopts an MDP-inspired formulation with iterative workspace reconstruction, where an evolving report serves as persistent memory and periodic synthesis mitigates context suffocation and noise contamination in long-horizon research.

Regarding state representation, approaches vary to ensure constant-size footprints. MEM1 (Zhou et al., 2025b) maintains a shared internal state that merges new observations with prior memory. Distinct from explicit text, MemGen (Zhang et al., 2025d) injects latent memory tokens directly into the reasoning stream.

**Hierarchical Folding**    For complex, long-horizon tasks, state maintenance requires structure beyond linear summarization. Hierarchical folding decomposes the task trajectory based on *subgoals*, maintaining fine-grained traces only while a subtask is active, and *folding* the completed sub-trajectory into a concise summary upon completion.

This *decompose-then-consolidate* strategy allows the working memory to expand and contract dynamically. HiAgent (Hu et al., 2025a) instantiates this by using subgoals as memory units, retaining only active action–observation pairs and writing back a summary after subgoal completion. Context-Folding (Zhang et al., 2025q) and AgentFold (Ye et al., 2025a) extend this by making the folding operation a learnable policy, training agents to autonomously determine when to branch into sub-trajectories and how to abstract them into high-level states. DeepAgent (Li et al., 2025h) further applies this to tool-use reasoning, compressing interactions into structured episodic and working memories to support fine-grained credit assignment. By replacing finished sub-trajectories with stable high-level abstractions, these methods preserve essential context while keeping the active window small.

**Cognitive Planning**    At the highest level of abstraction, working memory creates and maintains an externalized *plan* or *world model*. The state functions not merely as a summary of the past, but as a forward-looking structure that guides future actions.

PRIME (Tran et al., 2025) integrates retrieval directly into the planning loop, ensuring that memory updates actively support complex reasoning steps. In embodied and agentic environments, treating the language model as a high-level planner elevates the plan to the core of working memory. Approaches like SayPlan employ 3D scene graphs as queryable environmental memory to scale planning across large spaces (Rana et al., 2023). In GUI and household tasks, systems like Agent-S (Agashe et al., 2025) and KARMA (Wang et al., 2025q) stabilize long-horizon performance by anchoring reasoning to a hierarchical plan, using memory-augmented retrieval to bridge long-term knowledge with short-term execution.

By making plans and structured environment representations the readable and writable core of working

memory, agents can maintain goal consistency and revise strategies robustly against perception failures (Song et al., 2023).

**Summary**   Multi-turn working memory pivots on the construction of an operable **state carrier** rather than the retention of raw history. By integrating *state consolidation* to compress continuous streams, *hierarchical folding* to structure sub-trajectories, and *cognitive planning* to anchor future actions, these mechanisms effectively decouple reasoning performance from interaction length. This paradigm enables agents to maintain temporal coherence and goal alignment over indefinite horizons while adhering to strict computational and memory constraints.

## 5   Dynamics: How Memory Operates and Evolves?



**Figure 8** The operational dynamics of agent memory. We decouple the complete memory lifecycle into three fundamental processes that drive the system's adaptability and self-evolution: (1) **Memory Formation** transforms raw interactive experiences into information-dense knowledge units by selectively identifying patterns with long-term utility; (2) **Memory Evolution** dynamically integrates new memories into the existing repository through *consolidation*, *updating*, and *forgetting* mechanisms to ensure the knowledge base remains coherent and efficient; and (3) **Memory Retrieval** executes context-aware queries to access specific memory modules, thereby optimizing reasoning performance with precise information support. The alphabetical order denotes the sequence of operations within the memory systems.

The preceding sections introduced the architectural forms (Section 3) and functional roles (Section 4) of memory, outlining a relatively static conceptual framework for agent memory. However, such a static view

overlooks the inherent dynamism that fundamentally characterizes agentic memory. Unlike knowledge that is statically encoded in model parameters or fixed databases, an agentic memory system can dynamically construct and update its memory store and perform customized retrieval conditioned on different queries. This adaptive capability is crucial for enabling agents to self-evolve and engage in lifelong learning.

Accordingly, this section investigates the paradigm shift from static storage to dynamic memory management and utilization. This paradigm shift reflects the foundational operational advantages of agentic memory over static database approaches. In practice, an agentic memory system can autonomously extract refined, generalizable knowledge based on reasoning traces and environmental feedback. By dynamically fusing and updating this newly extracted knowledge with the existing memory base, the system ensures continuous adaptation to evolving environments and mitigates cognitive conflicts. Based on the constructed memory base, the system executes targeted retrieval from designated memory modules at precise moments, thereby enhancing reasoning effectively. To systematically analyze "how" the memory system operates and evolves, we examine the complete memory lifecycle by decomposing it into three fundamental processes. Figure 8 provides a holistic illustration of this dynamic memory lifecycle, highlighting how memory formation, evolution, and retrieval interact to support adaptive and self-evolving agent behavior.

> **Three Fundamental Process in Memory Systems**
>
> 1. **Memory Formation**(Section 5.1): This process focuses on transforming raw experience into information-dense knowledge. Instead of passively logging all interaction history, the memory system selectively identifies information with long-term utility, such as successful reasoning patterns or environmental constraints. This part answers the question: "How to extract the memory?".
>
> 2. **Memory Evolution**(Section 5.2): This process represents the dynamic evolution of the memory system. It focuses on integrating newly formed memories with the existing memory base. Through mechanisms such as the consolidation of correlated entries, conflict resolution, and adaptive pruning, the system ensures that the memory remains generalizable, coherent, and efficient in an ever-changing environment. This part answers the question: "How to refine the memory?".
>
> 3. **Memory Retrieval**(Section 5.3): This process determines the quality of the retrieved memory. Conditioned on the context, the system constructs a task-aware query and uses a carefully designed retrieval strategy to access the appropriate memory bank. The retrieved memory is therefore both semantically relevant and functionally critical for reasoning. This part answers the question: "How to utilize the memory?".

These three processes are not independent; rather, they form an interconnected cycle that drives the dynamic evolution and operation of the memory system. Memory extracted during the memory formation stage is integrated and updated with the existing memory base during the memory evolution stage. Leveraging the memory base constructed through these first two stages, the memory retrieval stage enables targeted access to optimize reasoning. In turn, reasoning outcomes and environmental feedback feed back into memory formation to extract new insights and into memory evolution to refine the memory base. Collectively, these components enable LLMs to transition from static conditional generators into dynamic systems that continuously learn from and respond to changing environments.

## 5.1 Memory Formation

We define memory formation as the process of encoding raw contexts (e.g., dialogues or images) into compact knowledge. The necessity for memory formation emerges from the scaling limitations inherent in processing lengthy, noisy, and highly redundant raw contexts. Full-context prompting often encounters computational overhead, prohibitive memory footprints, and degraded reasoning performance on out-of-distribution input lengths. To mitigate these issues, recent memory systems distill essential information into efficiently storable and precisely retrievable representations, enabling more efficient and effective inference.

Memory formation is not independent of the preceding sections. Depending on the task type, the memory formation process selectively extracts different architectural memories described in Section 3 to fulfill the

corresponding functions outlined in Section 4. Based on the granularity of information compression and the logic of encoding, we categorize the memory formation process into five distinct types. Table 7 summarizes representative methods under each category, comparing their sub-types, representation forms, and key mechanisms.

> **Five Categories of Memory Formation Operations**
>
> - **Semantic Summarization**(Section 5.1.1) transforms lengthy raw data into compact summaries, filtering out redundancy while preserving global, high-level semantic information to reduce contextual overhead.
>
> - **Knowledge Distillation** (Section 5.1.2) extracts specific cognitive assets, ranging from factual details to experiential planning strategies.
>
> - **Structured Construction**(Section 5.1.3) organizes amorphous source data into explicit topological representations, such as knowledge graphs or hierarchical trees, to enhance the explainability of memory and support multi-hop reasoning.
>
> - **Latent Representation**(Section 5.1.4) encodes raw experiences directly into machine-native formats (e.g., vector embeddings or KV states) within a continuous latent space.
>
> - **Parametric Internalization**(Section 5.1.5) consolidates external memories directly into the model's weight space through parameter updates, effectively transforming retrievable information into the agent's intrinsic competence and instincts.

Although we categorize these methods into five types, we argue that these memory formation strategies are not mutually exclusive. A single algorithm can integrate multiple memory formation strategies and translate knowledge across different representations (Li et al., 2025k).

### 5.1.1 Semantic Summarization

Semantic summarization transforms raw observational data into compact and semantically rich summaries. The resulting summaries capture the global, high-level information of the original data, rather than specific factual or experiential details (Zhao et al., 2024; Anokhin et al., 2024). Typical examples of such summaries include the overarching narrative of a document (Kim and Kim, 2025; Yu et al., 2025a), the procedural flow of a task (Ye et al., 2025a; Zhang et al., 2025q), or a user's historical profile (Zhang, 2024; Westhäußer et al., 2025). By filtering out redundant content while preserving task-relevant global semantics, semantic summarization provides a high-level guiding blueprint for subsequent reasoning without introducing excessive contextual overhead. To achieve these effects, the compression process can be implemented in two primary ways: incremental and partitioned semantic summarization.

**Incremental Semantic Summarization**   This paradigm employs a temporal integration mechanism that continuously fuses newly observed information with the existing summary, producing an evolving representation of global semantics. This chunk-by-chunk paradigm supports incremental learning (McCloskey and Cohen, 1989), circumvents the $O(n^2)$ computational burden of full-sequence processing (Yu et al., 2025a), and promotes progressive convergence toward global semantics (Chen et al., 2024b). Early implementations such as MemGPT (Packer et al., 2023a) and Mem0 (Chhikara et al., 2025) directly merged new chunks with existing summaries at appropriate moments, relying solely on the LLM's inherent summarization ability. However, this approach was constrained by the model's limited capacity, often resulting in inconsistency or semantic drift. To alleviate these issues, Chen et al. (2024b) and Wu et al. (2025i) incorporated external evaluators to filter redundant or incoherent content, using a convolutional-based discriminator for consistency verification and DeBERTa (He et al., 2020) for filtering trivial content, respectively. Instead of relying on auxiliary networks, subsequent methods, such as Mem1 (Zhou et al., 2025b) and MemAgent (Yu et al., 2025a), enhanced the LLM's own summarization capability through reinforcement learning with PPO (Schulman et al., 2017) and GRPO (Shao et al., 2024).

As incremental summarization advanced from heuristic fusion to filtered integration and ultimately to learning-based optimization, the summarization competence became increasingly internalized within the

**Table 7** Taxonomy of memory formation methods. We classify approaches based on the memory formation operations. Methods are analyzed across three technical dimensions: (1) **Sub-Type** identifies the specific variation or scope, (2) **Representation Form** specifies the output format, and (3) **Key Mechanism** denotes the core algorithmic strategy.

| Method | Sub-Type | Representation Form | Key Mechanism |
|---|---|---|---|
| *I. Semantic Summarization* | | | |
| MemGPT (Packer et al., 2023a) | Incremental | Textual Summary | Merging new chunks into the working context |
| Mem0 (Chhikara et al., 2025) | Incremental | Textual Summary | LLM-driven summarization |
| Mem1 (Zhou et al., 2025b) | Incremental | Textual Summary | RL-optimized summarization (PPO) |
| MemAgent (Yu et al., 2025a) | Incremental | Textual Summary | RL-optimized summarization (GRPO) |
| MemoryBank (Zhong et al., 2024) | Partitioned | Textual Summary | Daily/Session-based segmentation |
| ReadAgent (Lee et al., 2024a) | Partitioned | Textual Summary | Semantic clustering before summarization |
| LightMem (Fang et al., 2025b) | Partitioned | Textual Summary | Topic-clustered summarization |
| DeepSeek-OCR (Wei et al., 2025a) | Partitioned | Visual Token Mapping | Optical 2D mapping compression |
| FDVS (You et al., 2024) | Partitioned | Multimodal Summary | Multi-source signal integration (Subtitle/Object) |
| LangRepo (Kahatapitiya et al., 2025) | Partitioned | Multimodal Summary | Hierarchical video clip aggregation |
| *II. Knowledge Distillation* | | | |
| TiM (Liu et al., 2023a) | Factual | Textual Insight | Abstraction of dialogue into thoughts |
| RMM (Tan et al., 2025b) | Factual | Topic Insight | Abstraction of dialogue into topic-based memory |
| MemGuide (Du et al., 2025b) | Factual | User Intent | Capturing high-level user intent |
| M3-Agent (Long et al., 2025) | Factual | Text-addressable Facts | Compressing egocentric visual observations |
| AWM (Wang et al., 2024l) | Experiential | Workflow Patterns | Workflow extraction from success trajectories |
| Mem$^p$ (Fang et al., 2025d) | Experiential | Procedural Knowledge | Distilling gold trajectories into abstract procedures |
| ExpeL (Zhao et al., 2024) | Experiential | Experience Insight | Contrastive reflection and successful practices |
| R2D2 (Huang et al., 2025c) | Experiential | Reflective Insight | Reflection on reasoning traces vs. ground truth |
| $H^2R$ (Ye et al., 2025b) | Experiential | Hierarchical Insight | Two-tier reflection (Plan & Subgoal) |
| Memory-R1 (Yan et al., 2025b) | Experiential | Textual Knowledge | RL-trained LLMExtract module |
| Mem-$\alpha$ (Wang et al., 2025o) | Experiential | Textual Insight | Learnable insight extraction policy |
| *III. Structured Construction* | | | |
| KGT (Sun et al., 2024) | Entity-Level | User Graph | Encoding user preferences as nodes/edges |
| Mem0$^g$ (Chhikara et al., 2025) | Entity-Level | Knowledge Graph | LLM-based entity and triplet extraction |
| D-SMART (Lei et al., 2025) | Entity-Level | Dynamic Memory Graph | Constructing an OWL-compliant graph |
| GraphRAG (Edge et al., 2025) | Entity-Level | Hierarchical KG | Community detection and iterative summarization |
| AriGraph (Anokhin et al., 2024) | Entity-Level | Semantic+Episodic Graph | Dual-layer (Semantic nodes + Episodic links) |
| Zep (Rasmussen et al., 2025) | Entity-Level | Temporal KG | 3-layer graph (Episodic, Semantic, Community) |
| RAPTOR (Sarthi et al., 2024) | Chunk-Level | Tree Structure | Recursive GMM clustering and summarization |
| MemTree (Rezazadeh et al., 2025c) | Chunk-Level | Tree Structure | Bottom-up insertion and summary updates |
| H-MEM (Sun and Zeng, 2025) | Chunk-Level | Hierarchical JSON | Top-down 4-level hierarchy organization |
| A-MEM (Xu et al., 2025c) | Chunk-Level | Networked Notes | Discrete notes with semantic links |
| PREMem (Kim et al., 2025b) | Chunk-Level | Reasoning Patterns | Cross-session reasoning pattern clustering |
| CAM (Li et al., 2025f) | Chunk-Level | Hierarchical Graph | Disentangling overlapping clusters via replication |
| G-Memory (Zhang et al., 2025c) | Chunk-Level | Hierarchical Graph | 3-tier graph (interaction, query, insight) |
| *IV. Latent Representation* | | | |
| MemoryLLM (Wang et al., 2024j) | Textual | Latent Vector | Self-updatable latent embeddings |
| M+ (Wang et al., 2025m) | Textual | Latent Vector | Cross-layer long-term memory tokens |
| MemGen (Zhang et al., 2025d) | Textual | Latent Token | Latent memory trigger and weaver |
| ESR (Shen et al., 2024) | Multimodal | Latent Vector | Video-to-Language-to-Vector encoding |
| CoMEM (Wu et al., 2025d) | Multimodal | Continuous Embedding | Vision-language compression via Q-Former |
| Mem2Ego (Zhang et al., 2025l) | Multimodal | Multimodal Embedding | Embedding landmark semantics as latent memory |
| KARMA (Wang et al., 2025q) | Multimodal | Multimodal Embedding | Hybrid long/short-term memory encoding |
| *V. Parametric Internalization* | | | |
| MEND (Mitchell et al., 2022) | Knowledge | Gradient Decomposition | Auxiliary network for fast edits |
| ROME (Meng et al., 2022) | Knowledge | Model Parameters | Causal tracing and rank-one update |
| MEMIT (Meng et al., 2023) | Knowledge | Model Parameters | Mass-editing via residual distribution |
| CoLoR (Wistuba et al., 2023) | Knowledge | LoRA Parameters | Low-rank adapter training |
| ToolFormer (Schick et al., 2023) | Capability | Model Parameters | Supervised fine-tuning on API calls |

model, thereby reducing cumulative errors across iterations. Nevertheless, the serial update nature still poses computational bottlenecks (Fang et al., 2025b) and potential information forgetting, motivating the development of Partitioned semantic summarization approaches.

**Partitioned Semantic Summarization**  This paradigm adopts a spatial decomposition mechanism, dividing information into distinct semantic partitions and generating separate summaries for each. Early studies typically adopted heuristic partitioning strategies for handling long contexts. MemoryBank (Zhong et al., 2024) and COMEDY (Chen et al., 2025c) summarize and aggregate long-term dialogues by treating each day or session as a basic unit. Along the structural dimension, Wu et al. (2021) and Bailly et al. (2025) generate summaries of summaries by segmenting long documents into chapters or paragraphs. While intuitive, such approaches often suffer from semantic discontinuity across partitions. To address this issue, methods such as ReadAgent (Lee et al., 2024a) and LightMem (Fang et al., 2025b) introduce semantic or topical clustering before summarization, thereby enhancing inter-chunk coherence. Extending beyond textual compression, DeepSeek-OCR (Wei et al., 2025a) pioneers the idea of compressing long contexts via optical 2D mapping, achieving higher compression ratios in multimodal scenarios. In the video memory domain, FDVS (You et al., 2024) and LangRepo (Kahatapitiya et al., 2025) segment long videos into clips and generate textual summaries by integrating multi-source signals such as subtitles, object detection, and scene descriptions, which are then hierarchically aggregated into a global long video story.

Compared with incremental summarization, the partitioned approach offers superior efficiency and captures finer-grained semantics. However, its independent processing of each sub-chunk can lead to the loss of cross-partition semantic dependencies.

**Summary**  Semantic summarization operates as a lossy compression mechanism, aiming to distill the gist from lengthy interaction logs. Unlike verbatim storage, it prioritizes global semantic coherence over local factual precision, transforming linear streams into compact narrative blocks. The primary strength of semantic summarization is efficiency: it drastically reduces context length, making it ideal for long-term dialogue. However, the trade-off is resolution loss: specific details or subtle cues may be smoothed out, limiting their utility in evidence-critical tasks.

### 5.1.2 Knowledge Distillation

While semantic summarization captures the global semantics of raw data at a macro level, knowledge distillation operates at a finer granularity, extracting reusable knowledge from interaction trajectories or documents. In a broad sense, knowledge refers to the various forms of factual and experiential memory described in Section 4, depending on the task's underlying functions.

**Distilling Factual Memory**  This process focuses on transforming raw interactions and documents into explicit, declarative knowledge regarding users and environmental states. This process ensures that the agent maintains consistency and adaptability by retaining verifiable facts rather than transient context. In the domain of user modeling, systems such as TiM (Liu et al., 2023a) and RMM (Tan et al., 2025c) employ abstraction mechanisms to convert dialogue turns into high-level thoughts or topic-based memory, thereby preserving long-term persona coherence. For users' objective modeling, approaches like MemGuide (Du et al., 2025b) extract user intent descriptions from dialogues. During reasoning, it captures and updates goal states, separating confirmed constraints from unresolved intents to mitigate goal drift. Furthermore, this distillation extends to multimodal environments, where agents like ESR (Shen et al., 2024) and M3-Agent (Long et al., 2025) compress egocentric visual observations into text-addressable facts about object locations and user routines.

**Distilling Experiential Memory**  This process focuses on extracting the strategies underlying task execution from historical trajectories. By deriving planning principles from successful rollouts and corrective signals from failures, this paradigm enhances the agent's problem-solving ability on specific tasks. Through abstraction and generalization, it further supports cross-task knowledge transfer. As a result, experiential generalization enables the agent to continually refine its competence and move toward lifelong learning.

This line of research aims to derive high-level planning strategies and key insights from both successful and failed trajectories. Some approaches focus on success-based distillation, where systems such as AgentRR (Feng et al., 2025) and AWM (Wang et al., 2024l) summarize overall task plans from successful cases. Mem$^p$ (Fang et al., 2025d) analyzes and summarizes the gold trajectories from the training set, distilling them into abstract procedural knowledge. Others adopt failure-driven reflection, exemplified by Matrix (Liu et al., 2024), SAGE (Liang et al., 2025), and R2D2 (Huang et al., 2025c), which compare reasoning traces against ground-truth answers to identify error sources and extract reflective insights. Combining both, ExpeL (Zhao et al., 2024) and From Experience to Strategy (Xia et al., 2025) contrasts successful and failed experiences to uncover holistic planning insights.

However, prior work primarily focuses on summarizing task-level planning knowledge, lacking fine-grained, step-level insights. To address this gap, H$^2$R (Ye et al., 2025b) introduces a two-tier reflection mechanism: it follows ExpeL to construct a pool of high-level planning insights, while further segmenting trajectories by subgoal sequences to derive step-wise execution insights.

Earlier methods relied on fixed prompts for insight extraction, making their performance sensitive to prompt design and the underlying LLM's capacity. Recently, trainable distillation methods have become prevalent. Learn-to-Memorize (Zhang et al., 2025t) optimizes task-specific prompts for different agents. On the other hand, Memory-R1 (Yan et al., 2025b) uses an LLMExtract module to obtain experiential and factual knowledge, while only the subsequent fusion component is trained to integrate these outputs into the memory bank. Although these approaches adopt an end-to-end framework, they still fall short in enhancing the LLM's intrinsic ability to distill insights. To overcome this limitation, Mem-$\alpha$ (Wang et al., 2025o) explicitly trains the LLM on what insights to extract and how to preserve them.

**Summary**    This part focuses on extracting function-specific knowledge from the raw context, without addressing the structure of memory storage. Each piece of knowledge can be viewed as a flat memory unit. Simply storing multiple units in an unstructured table ignores the semantic and hierarchical relations among them. To address this, the memory formation process can apply structured rules to derive insights and store them within a hierarchical architecture. Simple but essential, the single knowledge distillation method introduced here serves as a foundational component for more complex and structured memory formation mechanisms.

### 5.1.3  Structured Construction

While Semantic Summarization (Section 5.1.1) and Knowledge Distillation (Section 5.1.2) effectively compress summaries and knowledge at different levels of granularity, they often treat memory as isolated units. In contrast, Structured Construction transforms amorphous data into organized topological representations. This process is not merely a change in storage format, but an active structural operation that determines how information is linked and layered. Unlike unstructured plaintext summarization, structured extraction significantly enhances both interpretability and retrieval efficiency. Crucially, such structural prior excels at capturing complex logic and dependencies in multi-hop reasoning tasks, offering substantial advantages over traditional retrieval-augmented methods.

Based on the operational granularity of how the underlying structure is derived, we categorize existing methods into two paradigms: Entity-Level Construction, which builds underlying topology by atomizing text into entities and relations, and Chunk-Level Construction, which builds structure by organizing intact text segments or memory items.

**Entity-Level Construction**    The foundational structure of this paradigm is derived from relational triples extraction, which decomposes the raw context into its finest-grained semantic atomic entities and relations. Traditional approaches model memory as a planar knowledge graph. For instance, KGT (Sun et al., 2024) introduces a real-time personalization mechanism where user preferences and feedback are directly encoded as nodes and edges within a user-specific knowledge graph. Similarly, Mem0$^g$ (Chhikara et al., 2025) utilizes LLMs to convert conversation messages directly into entities and relation triplets in the extraction phase.

However, these direct extraction methods are often limited by the inherent capabilities of the LLM, leading to potential noise or structural errors. To improve the quality of the constructed graph, D-SMART (Lei et al.,

2025) adopts a refined approach: it first employs an LLM to distill core semantic content into concise, assertion-like natural language statements, and subsequently extracts an OWL-compliant knowledge graph fragment through a neuro-symbolic pipeline. Additionally, Ret-LLM (Modarressi et al., 2023) applies supervised fine-tuning to the LLM, enabling more robust read-write interactions with the relational graph.

While the aforementioned methods focus on planar structures, recent advancements have progressed towards constructing hierarchical memory to capture high-level abstractions. For example, GraphRAG (Edge et al., 2025) derives an entity knowledge graph from source documents and applies community detection algorithms to extract graph communities and generate community summaries iteratively. This hierarchical approach identifies higher-level cluster associations between entities, enabling the extraction of generalized insights and facilitating flexible retrieval at varying granularities.

To better reflect the internal coherence and temporal information of the original data, some works extend the semantic knowledge graph by incorporating an episodic graph. AriGraph (Anokhin et al., 2024) and HippoRAG (Gutierrez et al., 2024) establish a dual-layer structure comprising the semantic and episodic graph. They extract semantic triplets from dialogues while connecting nodes that occur simultaneously or establishing node–paragraph indices. Zep (Rasmussen et al., 2025) further formalizes this into a three-layer temporal graph architecture: an episodic subgraph ($\mathcal{G}_e$) that logs the occurrence and processing times of raw messages via a bi-temporal model, a semantic subgraph ($\mathcal{G}_s$) for entities and time-bounded facts, and a community subgraph ($\mathcal{G}_c$) for high-level clustering and summarization of entities.

**Chunk-Level Construction**   This paradigm treats continuous text spans or discrete memory items as nodes, preserving local semantic integrity while organizing them into topological structures. The evolution of this field progresses from static, planar (2D) extraction from fixed corpora to dynamic adaptation with incoming trajectories, and ultimately to hierarchical (3D) architectures.

Early approaches focused on organizing fixed text libraries into static planar structures. HAT (A et al., 2024) processes long texts by segmenting them and progressively aggregating summaries to construct a hierarchical tree. Similarly, RAPTOR (Sarthi et al., 2024) recursively clusters text chunks using UMAP for dimensionality reduction and Gaussian Mixture Models for soft clustering, iteratively summarizing these clusters to form a tree. However, these static methods lack the flexibility to handle streaming data without costly reconstruction.

To address this, dynamic planar approaches incrementally build memory structures as new trajectories arrive, differing based on their foundational elements. Methods based on raw text include MemTree (Rezazadeh et al., 2025c) and H-MEM (Sun and Zeng, 2025). MemTree adopts a bottom-up approach where new text fragments retrieve the most similar nodes and are inserted as children or iteratively into a subtree, triggering bottom-up summary updates for all parent nodes. Conversely, H-MEM utilizes a top-down strategy, prompting the LLM to organize data into a four-level JSON hierarchy comprising the domain, category, memory trace, and episode layer. Alternatively, A-MEM (Xu et al., 2025c) and PREMem (Kim et al., 2025b) focus on reorganizing the extracted memory items. A-MEM summarizes knowledge into discrete notes and links relevant ones to construct a networked memory. PREMem clusters extracted factual, experiential, and subjective memories to identify and store higher-dimensional cross-session reasoning patterns.

Recent advancements move beyond planar layouts to construct hierarchical structures, offering richer semantic depth. SGMem (Wu et al., 2025h) constructs a hierarchy by using NLTK to split text into sentences, forming a KNN graph across all sentence nodes, and subsequently calling an LLM to extract summaries, facts, and insights corresponding to each dialogue. To support the incremental construction of hierarchical structures as streaming data arrives, CAM (Li et al., 2025f) establishes edges between text blocks based on semantic relevance and narrative coherence. It iteratively summarizes the ego graph and handles new memory insertions by explicitly disentangling overlapping clusters through node replication. In multi-agent scenarios, G-memory (Zhang et al., 2025c) extends this dynamic 3D approach by maintaining three distinct graphs: an interaction graph for raw chat history, a query graph for specific tasks, and an insight graph. This structure enables each agent to receive customized memory at varying levels of granularity.

**Summary**   The main advantage of structured construction is explainability and the ability to handle complex relational queries. Such methods capture intricate semantic and hierarchical relationships between memory elements, support reasoning over multi-step dependencies, and facilitate integration with symbolic or

graph-based reasoning frameworks. However, the downside is schema rigidity: pre-defined structures may fail to represent nuanced or ambiguous information, and the extraction and maintenance costs are typically high.

### 5.1.4 Latent Representation

The previous chapters focused on how to build token-level memory; this part focuses on encoding memory into the machine's native latent representation. Latent representation encodes raw experiences into embeddings that reside in a latent space. Unlike semantic compression and structured extraction, which summarize experiences before embedding them into vectors, latent encoding inherently stores experiences in latent space, thereby reducing information loss during summarization and text embedding. Furthermore, latent encoding is more conducive to machine cognition, enabling a unified representation across different modalities and ensuring both density and semantic richness in memory representation.

**Textual Latent Representation**   Although originally designed to accelerate inference, the KV cache can also be viewed as a form of latent representation within the context of memory (Li et al., 2025c; Jiang et al., 2025b). It utilizes additional memory to store past information, thereby avoiding redundant computation. MEMORYLLM (Wang et al., 2024j) and M+ (Wang et al., 2025m) represent memory as self-updatable latent embeddings, which are injected into transformer layers during inference. Moreover, MemGen (Zhang et al., 2025d) introduces a memory trigger to monitor the agent's reasoning state and determine when to explicitly invoke memory, as well as a memory waiver that leverages the agent's current state to construct a latent token sequence. This sequence acts as machine-native memory, enriching the agent's reasoning capabilities.

**Multimodal Latent Representation**   In multimodal memory research, CoMEM (Wu et al., 2025d) compresses vision-language inputs into fixed-length tokens via a Q-Former, enabling dense, continuous memory and supporting plug-and-play usage for infinite context lengths. Encode-Store-Retrieve (Shen et al., 2024) converts egocentric video frames into language encodings using Ego-LLaVA, which are subsequently transformed into vector representations through an embedding model. Although embedding models are employed to ensure semantic alignment, these methods often face a trade-off between compression loss and computational overhead, particularly in handling gradient flow in long-context sequences.

When integrated with Embodied AI, multimodal latent memory can fuse data from multiple sensors. For example, Mem2Ego (Zhang et al., 2025l) dynamically aligns global contextual information with local perception, embedding landmark semantics as latent memory to enhance spatial reasoning and decision-making in long-horizon tasks. KARMA (Wang et al., 2025q) adopts a hybrid long- and short-term memory form that encodes object information into multimodal embeddings, achieving a balance between immediate responsiveness and consistent representation. These explorations underscore the advantages of latent encoding in providing unified and semantically rich representations across modalities.

**Summary**   Latent representation bypasses human-readable formats, encoding experiences directly into machine-native vectors or KV-caches. This high-density format preserves rich semantic signals that might be lost in text decoding, enabling smoother integration with the model's internal computations. And it supports multimodal alignment seamlessly. However, it suffers from opaqueness. The latent memory is a black box, making it difficult for humans to debug, edit, or verify the knowledge it stores.

### 5.1.5 Parametric Internalization

As LLMs increasingly incorporate memory systems to support long-term adaptation, a central research question is how these external memories should be consolidated into parametric form. While the latent representation methods discussed above parameterize memory externally to the model, parametric internalization directly adjusts the model's internal parameters. It leverages the model's capacity to encode and generalize information through its learned parameter space. This paradigm fundamentally enhances the model's intrinsic competence, eliminating the overhead of external storage and retrieval while seamlessly supporting continual updates. As we discussed in Section 4, not all memory content serves the same function: some entries provide declarative knowledge, while others encode procedural strategies that shape an agent's reasoning and behavior. This distinction motivates a finer-grained view of memory internalization, separating it into knowledge internalization and capability internalization.

**Knowledge Internalization**  This strategy entails converting externally stored factual memories, such as conceptual definitions or domain knowledge, into the model's parameter space. Through this process, the model can directly recall and utilize these facts without relying on explicit retrieval or external memory modules. In practice, knowledge internalization is typically achieved through model editing (Sinitsin et al., 2020; De Cao et al., 2021). Early work, such as MEND (Mitchell et al., 2022), introduced an auxiliary network that enables rapid, single-step edits by decomposing fine-tuning gradients, thereby minimizing interference with unrelated knowledge. Building on this line of work, ROME (Meng et al., 2022) refined the editing process by using causal tracing to precisely locate the MLP layers that store specific facts and applying rank-one updates to inject new information with higher precision and better generalization. MEMIT (Meng et al., 2023) further advanced this line by supporting batch edits, enabling thousands of facts to be updated simultaneously through multi-layer residual distributions and batch formulas, which substantially improves scalability. With the rise of parameter-efficient paradigms like LoRA (Hu et al., 2022), knowledge internalization can be performed through lightweight adapters rather than direct parameter modification. For instance, CoLoR (Wistuba et al., 2023) freezes the pretrained Transformer parameters and trains only small LoRA adapters to internalize new knowledge, avoiding the high cost of full-parameter fine-tuning. Despite these advances, these approaches can still incur off-target effects (De Cao et al., 2021) and remain vulnerable to catastrophic forgetting in continual learning scenarios.

**Capability Internalization**  This strategy seeks to embed experiential knowledge, such as procedural expertise or strategic heutistics, into the model's parameter space. The paradigm represents a memory formation operation in a broad sense, shifting from the acquisition of factual knowledge to the internalization of experiential capabilities. Specifically, these capabilities include domain-specific solution schemas, strategic planning, and the effective deployment of agentic skills, among others. Technically, capability internalization is achieved by learning from reasoning traces, through supervised fine-tuning (Wei et al., 2022; Zelikman et al., 2022; Schick et al., 2023; Mukherjee et al., 2023) or preference-guided optimization methods such as DPO (Rafailov et al., 2023; Tunstall et al., 2023; Yuan et al., 2024c; Grattafiori et al., 2024) and GRPO (Shao et al., 2024; DeepSeek-AI et al., 2025). Since this aspect does not fall within the scope of typical agentic memory research, it will not be discussed in detail in this section.

**Summary**  Parametric internalization represents the ultimate consolidation of memory, where external knowledge is fused into the model's weights via gradients. This shifts the paradigm from retrieving information to possessing capability, mimicking biological long-term potentiation. As knowledge becomes effectively instinctive, access is zero-latency, enabling the model to respond immediately without querying external memory. However, this approach faces several challenges, including catastrophic forgetting and high update costs. Unlike external memory, parameterized internalization is difficult to modify or remove precisely without unintended side effects, limiting flexibility and adaptability.

## 5.2  Memory Evolution

Memory Formation introduced in Section 5.1 extracts memory from raw data. The next important step is to integrate the newly extracted memories with the existing memory repository, enabling the dynamic evolution of the memory system. A naive strategy is simply appending new entries to the existing memory bank. However, it overlooks the semantic dependencies and potential contradictions between memory entries and neglects the temporal validity of information. To address these limitations, we introduce Memory Evolution. This mechanism consolidates new and existing memories to synthesize high-level insights, resolve logical conflicts, and prune obsolete data. By ensuring the compactness, consistency, and relevance of long-term knowledge, this approach enables the memory system to adapt its cognitive processes and contextual understanding as environments and tasks evolve.

Based on the objectives of memory evolution, we categorize it into the following mechanisms:

> **Three Mechanisms of Memory Evolution**
>
> - **Memory Consolidation** (Section 5.2.1) merges new and existing memories and performs reflective integration, forming more generalized insights. This ensures that learning is cumulative rather than isolated.
>
> - **Memory Updating** (Section 5.2.2) resolves conflicts between new and existing memories, correcting and supplementing the repository to maintain accuracy and relevance. It allows the agent to adapt to changes in the environment or task requirements.
>
> - **Memory Forgetting** (Section 5.2.3) removes outdated or redundant information, freeing capacity and improving efficiency. This prevents performance degradation due to knowledge overload and ensures that the memory repository remains focused on actionable and current knowledge.

These mechanisms collectively maintain the generalization, accuracy, and timeliness of the memory repository. By actively managing memory evolution, these mechanisms underscore the agentic capabilities of the memory system, facilitating continuous learning and autonomous self-improvement. Figure 9 provides a unified view of these memory evolution mechanisms, illustrating their operational roles and representative frameworks within a shared memory database.

### 5.2.1 Consolidation

Memory consolidation aims to transform newly acquired short-term traces into structured and generalizable long-term knowledge. Its core mechanism is to identify semantic relationships between new and existing memories and to integrate them into higher-level abstractions or insights. This process serves two main purposes. First, it reorganizes fragmented pieces of information into coherent structures, preventing the loss of critical details during short-term retention and enabling the formation of stable knowledge schemas. Second, by abstracting, compressing, and generalizing experiential data, consolidation extracts reusable patterns from specific events, yielding insights that support cross-task generalization.

A central challenge is determining the granularity at which new memories should be matched and merged with existing ones. Prior work spans a spectrum of consolidation strategies, from local content merging to cluster-level fusion and global integration.

**Local Consolidation**   This operation focuses on fine-grained updates involving highly similar memory fragments. In RMM (Tan et al., 2025c), each new topic memory retrieves its top-K most similar candidates, and an LLM decides whether merging is appropriate, thereby reducing the risk of incorrect generalization. In multimodal settings, VLN (Song et al., 2025b) triggers a pooling mechanism when capacity is saturated. It identifies the most similar or redundant memory pairs and compresses them into higher-level abstractions. These approaches refine detailed knowledge while preserving the global structure of the memory store, improving precision and storage efficiency. However, they cannot fully capture cluster-level relations or the higher-order dependencies that emerge across semantically related memories.

**Cluster-level Fusion**   Adopting cluster-level fusion is essential for capturing cross-instance regularities as memory grows. Across clusters, PREMem (Kim et al., 2025b) aligns new memory clusters with similar existing ones and applies fusion modes such as generalization and refinement to form higher-order reasoning units, substantially improving interpretability and reasoning depth. Within a cluster, TiM (Liu et al., 2023a) periodically invokes an LLM to examine memories that share the same hashing bucket and merges semantically redundant entries. CAM (Li et al., 2025f) merges all nodes within the target cluster into a representative summary, yielding higher-level and consistent cross-sample representations. These methods reorganize the memory structure at a broader scale and mark an important step toward structured knowledge.

**Global Integration**   This operation performs holistic consolidation to maintain global coherence and to distill system-level insights from accumulated experience. Compared with Section 5.1.1, semantic summarization focuses on deriving a global summary from the existing context and can be viewed as the initial construction of the summary. In contrast, this paragraph emphasizes how new information is integrated into an existing

**Figure 9** The landscape of Memory Evolution mechanisms. We categorize the evolution process into three distinct branches that maintain the central **Memory Database**: (a) **Consolidation** synthesizes insights by processing raw materials through local consolidation, cluster fusion, and global integration; (b) **Updating** ensures accuracy and consistency by performing conflict resolution on external databases and applying parameter updates to the internal model; and (c) **Forgetting** optimizes efficiency by pruning data based on specific criteria: time expiration, low access frequency, and low informational value. The outer ring displays representative frameworks and agents associated with each evolutionary mechanism.

summary as additional data arrives. For user factual memory, MOOM (Chen et al., 2025d) constructs stable role profiles by integrating temporary role snapshots with historical traces using rule-based processing, embedding methods, and LLM-driven abstraction. For experiential memory, Matrix (Liu et al., 2024) performs iterative optimization to combine execution trajectories and reflective insights with global memory, distilling task-agnostic principles that support reuse across scenarios. As single-step reasoning contexts and environmental feedback lengthen, methods like AgentFold (Ye et al., 2025a) and Context Folding (Zhang et al., 2025q) internalize the ability to compress working memory. In multi-step interactions, including web navigation, these methods automatically summarize and condense the global context after each step, supporting efficient and effective reasoning. Global integration consolidates high-level, structured knowledge from the complete history of experience, providing a reliable contextual foundation while improving generalization, reasoning accuracy, and personalized decision-making.

**Summary** Consolidation is the cognitive process of reorganizing fragmented short-term traces into coherent long-term schemas. It moves beyond simple storage to synthesize connections between isolated entries,

forming a structured worldview. It enhances generalization and reduces storage redundancy. However, it risks information smoothing, where outlier events or unique exceptions are lost during the abstraction process, potentially reducing the agent's sensitivity to anomalies and specific events.

### 5.2.2 Updating

Memory Update refers to the process by which an agent revises or replaces its existing memory when conflicts arise or new information is acquired. The goal is to maintain factual consistency and continual adaptation without full model retraining. Unlike memory consolidation described in Section 5.2.1, which focuses on abstraction and generalization, memory update emphasizes localized correction and synchronization, enabling the agent to remain aligned with an evolving environment.

Through continuous updating, agentic memory systems preserve the accuracy and timeliness of knowledge, preventing outdated information from biasing reasoning. It is thus a core mechanism for achieving lifelong learning and self-evolution. Depending on where the memory resides, updates fall into two categories: (1) External Memory Update: updates to external memory stores and (2) Model Editing: model-internal editing within the parameter space.

**External Memory Update**  Entries in vector databases or knowledge graphs are revised whenever contradictions or new facts emerge. Instead of altering model weights, this approach maintains factual alignment through dynamic modifications of external storage. Static memories inevitably accumulate outdated or conflicting entries, leading to logical inconsistencies and reasoning errors. Updating external memories enables lightweight corrections while avoiding the cost of full retraining or re-indexing.

The development of external memory update mechanisms has progressed along a trajectory, moving from rule-based corrections to temporally aware soft deletion, then to delayed-consistency strategies, and ultimately to fully learned update policies. Early systems such as MemGPT (Packer et al., 2023a), D-SMART (Lei et al., 2025), and Mem0$^g$ (Chhikara et al., 2025) followed a straightforward pipeline in which the LLM detects conflicts between new information and then invokes replace or delete operations to update the memory. Although effective for basic factual repair, these systems relied on destructive replacement, erasing valuable historical context and breaking temporal continuity. To address this issue, Zep (Rasmussen et al., 2025) introduced temporal annotations, marking conflicting facts with invalid timestamps rather than deleting them, thereby preserving both semantic consistency and temporal integrity. This marked a shift from hard replacement to soft, time-aware updating. However, real-time updates impose significant computational and I/O burdens under high-frequency interaction. MOOM (Chen et al., 2025d) and LightMem (Fang et al., 2025b), therefore, introduced dual-phase updating: a soft online update for real-time responsiveness, followed by an offline reflective consolidation phase where similar entries are merged and conflicts resolved via LLM reasoning. This eventual consistency paradigm balances latency and coherence. As agentic reinforcement learning matured, it became possible to enhance the LLM's intrinsic memory update decision-making through reinforcement learning. Mem-$\alpha$ (Wang et al., 2025o) formulated memory updating as a policy-learning problem, enabling the LLM to learn when, how, and whether to update, thereby achieving dynamic trade-offs between stability and freshness.

Overall, external memory updates have transitioned from manually triggered corrections to self-regulated, temporally aware learning processes, maintaining factual consistency and structural stability through LLM-driven retrieval, conflict detection, and revision.

**Model Editing**  Model editing performs direct modifications within the model's parameter space to correct or inject knowledge without full retraining, representing implicit knowledge updates. Retraining is costly and prone to catastrophic forgetting. Model editing enables precise, low-cost corrections that enhance adaptability and internal knowledge retention.

Approaches of model editing fall into two main categories. (1) Explicit localization and modification: ROME (Tan et al., 2025b) identifies the parameter region encoding specific knowledge via gradient tracing and performs targeted weight updates; Model Editor Networks (Tang et al., 2025c) trains an auxiliary meta-editor network to predict optimal parameter adjustments. (2) Latent-space self-updating: MEMORYLLM (Xu et al., 2025c) embeds a memory pool within Transformer layers, periodically replacing memory tokens to integrate

new knowledge; M+ (Wang et al., 2025m) maintains dual-layer memories, discarding obsolete short-term entries and compressing key information into long-term storage.

Hybrid approaches such as ChemAgent (Tang et al., 2025c) further combine external memory updates with internal model editing, synchronizing factual and representational changes for rapid cross-domain adaptation.

**Summary**   From an implementation standpoint, memory updating focuses on resolving conflicts and revising knowledge triggered by the arrival of new memories, whereas memory consolidation emphasizes the integration and abstraction of new and existing knowledge. The two memory updating strategies discussed above establish a dual-pathway mechanism involving conflict resolution in external databases and parameter editing within the model, enabling agents to perform continuous self-correction and support long-term evolution. The key challenge is the stability–plasticity dilemma: determining when to overwrite existing knowledge versus when to treat new information as noise. Incorrect updates can overwrite critical information, leading to knowledge degradation and faulty reasoning.

### 5.2.3   Forgetting

Memory forgetting refers to the deliberate removal of outdated, redundant, or low-value information to free capacity and maintain focus on salient knowledge. Unlike update mechanisms, which resolve conflicts between memories, forgetting prioritizes eliminating outdated information to ensure efficiency and relevance. Over time, unbounded memory accumulation leads to increased noise, retrieval delays, and interference from outdated knowledge. Controlled forgetting helps mitigate overload and maintain cognitive focus. Yet, overly aggressive pruning risks erasing rare but essential knowledge, harming reasoning continuity in long-term contexts.

Forgetting mechanisms can be categorized into Time-based Forgetting, Frequency-based Forgetting, and Importance-driven Forgetting, corresponding respectively to creation time, retrieval activity, and integrated semantic valuation.

**Time-based Forgetting**   Time-driven forgetting considers only the creation time of memories, gradually decaying their strength over time to emulate human memory fading. MemGPT (Packer et al., 2023a) evicts the earliest messages upon context overflow. Xu et al. (2025c) and Wang et al. (2025m) employ stochastic token replacement, with a replacement ratio of K/N, to simulate exponential forgetting in human cognition, discarding the oldest entries once the pool exceeds capacity. Unlike explicit deletion of old memories, MAICC (Jiang et al., 2025c) implements soft forgetting by gradually decaying the weights of memories over time. This process mirrors natural forgetting, ensuring continuous adaptation without historical overload.

**Frequency-based Forgetting**   Frequency-driven forgetting prioritizes memory based on retrieval behavior, retaining frequently accessed entries while discarding inactive ones. XMem (Cheng and Schwing, 2022) employs an LFU policy to remove low-frequency entries; KARMA (Wang et al., 2025q) uses counting Bloom filters to track access frequency; MemOS (Li et al., 2025k) applies an LRU strategy, removing long-unused items while archiving highly active ones. This ensures efficient retrieval and storage equilibrium. By distinguishing between creation time and retrieval frequency, these two axes form a more orthogonal taxonomy: time-based decay captures natural temporal aging, while frequency-based forgetting reflects usage dynamics, together maintaining system efficiency and recency.

**Importance-driven Forgetting**   Importance-driven forgetting integrates temporal, frequency, and semantic signals to retain high-value knowledge while pruning redundancy. Early works such as Zhong et al. (2024) and Chen et al. (2025d) quantified importance via composite scores combining temporal decay and access frequency, achieving numeric-based selective forgetting. Later methods evolved toward semantic-level evaluation: VLN (Song et al., 2025b) pools semantically redundant memories via similarity clustering, while Livia (Xi and Wang, 2025) incorporates emotional salience and contextual relevance to model emotion-driven selective forgetting. As LLMs develop increasingly powerful judgment capabilities, TiM (Liu et al., 2023a) and MemTool (Lumer et al., 2025) leverage LLMs to assess memory importance and explicitly prune or forget less important memories. This shift reflects a transition from static numeric scoring to semantic intelligence. Agents can now perform conscious forgetting and selectively retain memories most pertinent to the task context, semantics, and affective cues.

**Summary** Time-based decay reflects the natural temporal fading of memory, frequency-based forgetting ensures efficient access to frequently used memories, and importance-driven forgetting introduces semantic discernment. These three forgetting mechanisms jointly govern how agentic memory remains timely, efficiently accessible, and semantically relevant. However, heuristic forgetting mechanisms like LRU may eliminate long-tail knowledge, which is seldom accessed but essential for correct decision-making. Therefore, when storage cost is not a critical constraint, many memory systems avoid directly deleting certain memories.

## 5.3 Memory Retrieval

Building on the memory bank established in Section 5.1 and Section 5.2, the next critical step is how to retrieve and utilize memories during reasoning. We define memory retrieval as the process of retrieving relevant and concise knowledge fragments from a certain memory repository to support current reasoning tasks at the right moment. The key challenge lies in efficiently and accurately locating the required knowledge fragments within a large-scale memory store. To address this, many algorithms employ heuristic strategies or learnable models to optimize various stages of the retrieval process. Based on the execution order of retrieval, this process can be decomposed into four aspects. Figure 10 provides a structured overview of this retrieval pipeline, organizing existing methods according to their roles across retrieval stages.

**Four Steps of Memory Retrieval**

- **Retrieval Timing and Intent**(Section 5.3.1) determines the specific moments and objectives for memory retrieval, shifting from passive, instruction-driven triggers to autonomous, self-regulated decisions.

- **Query Construction**(Section 5.3.2) bridges the semantic gap between the user's raw input and the stored memory index by decomposing or rewriting queries into effective retrieval signals.

- **Retrieval Strategies**(Section 5.3.3) executes the search over the memory repository, employing paradigms ranging from sparse lexical matching to dense semantic embedding and structure-aware graph traversal.

- **Post-Retrieval Processing**(Section 5.3.4) refines the retrieved raw fragments through re-ranking, filtering, and aggregation, ensuring that the final context provided to the model is concise and coherent.

Collectively, these mechanisms transform memory retrieval from a static search operation into a dynamic cognitive process. Retrieval timing and intent determine when and where to retrieve. Next, query construction specifies what to retrieve, and retrieval strategies focus on how to execute the retrieval. Finally, post-retrieval processing decides how the retrieved information is integrated and used. A robust agentic system typically orchestrates these components within a unified pipeline, enabling agents to approximate human-like associative memory activation for efficient knowledge access.

### 5.3.1 Retrieval Timing and Intent

The retrieval intent and timing determine when to trigger the retrieval mechanism and which memory store to query. Existing memory systems adopt different design choices in this regard, ranging from always-on retrieval to retrieval triggered by explicit instructions or internal signals (Zhao et al., 2024; Wang et al., 2025o; Fang et al., 2025b). For example, MIRIX (Wang and Chen, 2025) performs retrieval from all six memory databases for each query and concatenates the retrieved contents, reflecting a design that prioritizes comprehensive memory access. Other approaches instead aim to trigger retrieval more selectively, allowing the model to decide both the timing and scope of memory access, which can lead to more targeted and efficient use of memory resources. In this subsection, we review the literature from two complementary perspectives: automated retrieval timing and automated retrieval intent.

**Automated Retrieval Timing** This term refers to the model's ability to autonomously determine when to trigger a memory retrieval operation during reasoning. The simplest strategy is to delegate the decision to either the LLM or an external controller, allowing it to determine solely from the query whether retrieval is necessary. For example, MemGPT (Packer et al., 2023a) and MemTool (Lumer et al., 2025) allow the LLM

**Figure 10** Taxonomy of memory retrieval methodologies in agentic systems. The mindmap organizes existing literature into four distinct phases of the retrieval pipeline: **Timing and Intent**, which governs the initiation of the process; **Query Construction**, covering techniques for query decomposition and rewriting; **Retrieval Strategies**, categorizing search paradigms into lexical, semantic, graph-based, and hybrid approaches; and **Post-Retrieval Processing**, which focuses on refining outputs through re-ranking, filtering, and aggregation.

itself to invoke retrieval functions, enabling efficient access to external memory within an operating-system-like framework. However, these methods rely on static judgments from the query alone, neglecting the model's dynamically evolving cognitive state during reasoning.

To address this limitation, recent work integrates fast–slow thinking mechanisms into retrieval timing. ComoRAG (Wang et al., 2025f) and PRIME (Tran et al., 2025), for instance, first produce a fast response and then let the agent evaluate its adequacy. If the initial reasoning is deemed insufficient, the system triggers deeper retrieval and reasoning based on failure feedback. MemGen (Zhang et al., 2025d) further refines the triggering mechanism by converting the explicit agent-level decision into a latent, trainable process. It introduces memory triggers that detect critical retrieval moments from latent rollout states, thereby improving the precision of retrieval timing while preserving end-to-end differentiability.

**Automated Retrieval Intent**   This aspect concerns the model's ability to autonomously decide which memory source to access within a hierarchical storage form. AgentRR (Feng et al., 2025), for example, dynamically switches between low-level procedural templates and high-level experiential abstractions based on environmental feedback. However, its reliance on explicit feedback limits applicability in open-ended reasoning settings.

To overcome this constraint, MemOS (Li et al., 2025k) employs a MemScheduler that dynamically selects among parametric, activation, and plaintext memory based on user-, task-, or organization-level context. Yet, this flat selection scheme overlooks the hierarchical structure of the memory system. H-MEM (Sun and Zeng, 2025) addresses this by introducing an index-based routing mechanism, which performs coarse-to-fine retrieval, moving from the domain layer to the episode layer and gradually narrowing the search space to the most relevant sub-memories. This hierarchical routing not only improves retrieval precision but also mitigates information overload.

**Summary** Autonomous timing and intent help reduce computational overhead and suppress unnecessary noise, but they also create a potential vulnerability. When an agent overestimates its internal knowledge and fails to initiate retrieval when needed, the system can fall into a silent failure mode in which knowledge gaps may lead to hallucinated outputs. Therefore, a balance needs to be achieved: providing the agent with essential information at the right moments while avoiding excessive retrieval that introduces noise.

### 5.3.2 Query Construction

After initiating the retrieval process, the next challenge lies in transforming the raw query into an effective retrieval signal aligned with the memory index. Query construction acts as the translation layer between the user's surface utterance and the memory's latent storage. Traditional approaches typically perform retrieval directly based on the user query, which is simple but fails to align the query semantics with those of the memory index. To bridge this gap, agentic memory systems proactively perform query decomposition or query rewriting, generating intermediate retrieval signals that better match the latent structure of the memory.

**Query Decomposition** This approach breaks down a complex query into simpler sub-queries, allowing the system to retrieve more fine-grained and relevant information. Such decomposition alleviates the one-shot retrieval bottleneck by enabling modular retrieval and reasoning over intermediate results. For instance, Visconde (Pereira et al., 2023) and ChemAgent (Tang et al., 2025c) employ LLMs to decompose the original question into sub-problems, retrieve candidate results for each from the memory, and finally aggregate them into a coherent answer. However, these methods lack global planning. To address this issue, PRIME (Tran et al., 2025) and MA-RAG (Nguyen et al., 2025) introduce a Planner Agent, inspired by the ReAct (Yao et al., 2023b) paradigm, that first formulates a global retrieval plan before decomposing it into sub-queries. Yet, these approaches mainly rely on problem-driven decomposition and thus cannot explicitly identify what specific knowledge the model is missing. To make sub-queries more targeted, Agent KB (Tang et al., 2025d) adopts a two-stage retrieval process in which a teacher model observes the student model's failures and generates fine-grained sub-queries accordingly. This targeted decomposition improves retrieval precision and reduces irrelevant results, particularly in knowledge-intensive tasks.

**Query Rewriting** Instead of decomposing, this strategy rewrites the original query or generates a hypothetical document to refine its semantics before retrieval. Such rewriting mitigates the mismatch between user intent and the memory index. HyDE (Gao et al., 2023b), for example, instructs the LLM to generate a hypothetical document in a zero-shot manner and performs retrieval using its semantic embedding. The generated document encapsulates the desired semantics, effectively bridging the gap between the user query and the target memory. MemoRAG (Qian et al., 2025) extends this idea by incorporating global memory into hypothetical document generation. It first compresses the global memory and then generates a draft answer conditioned on both the query and the compressed memory; this draft is then used as a rewritten query. Since the draft has access to the global memory context, it captures user intent more faithfully and uncovers implicit information needs. Similarly, MemGuide (Du et al., 2025b) leverages the dialogue context to prompt an LLM to produce a concise, command-like phrase that serves as a high-level intent description for retrieval. Beyond directly prompting an LLM to rewrite the query, Rewrite-Retrieve-Read (Ma et al., 2023b) trains a small language model as a dedicated rewriter through reinforcement learning, while ToC (Kim et al., 2023a) employs a Tree of Clarifications to progressively refine and specify the user's retrieval objective.

**Summary** These two paradigms, decomposition and rewriting, are not mutually exclusive. Auto-RAG (Kim et al., 2024a) integrates both by evaluating HyDE and Visconde under identical retrieval conditions and then

selecting the strategy that performs best for the given task. The findings of this work demonstrate that the quality of the memory-retrieval query has a substantial impact on reasoning performance. In contrast to earlier research, which primarily focused on designing sophisticated memory architectures, recent studies (Yan et al., 2025a) place increasing emphasis on the retrieval construction process, shifting the role of memory toward serving retrieval. The choice of what to retrieve with is, unsurprisingly, a critical component of this process.

### 5.3.3  Retrieval Strategies

After clarifying the retrieval objective, we obtain a query with a well-defined intent. The next core challenge lies in leveraging this query to efficiently and accurately retrieve truly relevant knowledge from a large and complex memory repository. Retrieval strategies serve as the bridge between queries and the memory base, and their design directly determines both retrieval efficiency and result quality. In this section, we systematically review various retrieval paradigms and analyze their strengths, limitations, and application scenarios—from traditional sparse retrieval based on keyword matching, to modern dense retrieval using semantic embeddings, to graph-based retrieval for structured knowledge, to the emerging class of generative retrieval methods, and finally to hybrid retrieval techniques that integrate multiple paradigms.

**Lexical Retrieval**  This strategy relies on keyword matching to locate relevant documents, with representative methods including TF-IDF (SPARCK JONES, 1972) and BM25 (Robertson and Zaragoza, 2009). TF-IDF measures the importance of keywords based on term frequency and inverse document frequency, enabling fast and interpretable retrieval. BM25 further refines this approach by incorporating term frequency saturation and document length normalization. Such methods are often employed in precision-oriented retrieval scenarios, where accuracy and relevance of results take precedence over recall (Tang et al., 2025d; Wang et al., 2025o; Pan et al., 2025). However, purely lexical matching struggles to capture semantic variations and contextual relationships, making it highly sensitive to linguistic expression differences and thus less effective in open-domain knowledge or multimodal memory settings.

**Semantic Retrieval**  This strategy encodes queries and memory entries into a shared embedding space and matches them based on semantic similarity rather than lexical overlap. Representative approaches utilize semantic encoders, including Sentence-BERT (Reimers and Gurevych, 2019) and CLIP (Radford et al., 2021). Within memory systems, this approach better captures task context and supports semantic generalization and fuzzy matching, making it the default choice in most agentic memory frameworks (Lewis et al., 2020; Wang et al., 2024b; Yang et al., 2024a; Xu et al., 2025c; Tan et al., 2025c; Nguyen et al., 2025; Qian et al., 2025; Hassell et al., 2025; Huang et al., 2025c). However, semantic drift and forced top-K retrieval often introduce retrieval noise and spurious recall. To address these issues, recent systems incorporate dynamic retrieval policies, reranking modules, and hybrid retrieval schemes.

**Graph Retrieval**  This strategy leverages not only semantic signals but also the explicit topological structure of graphs, enabling inherently more precise and structure-aware retrieval. By directly accessing structural paths, these methods exhibit stronger multi-hop reasoning capabilities and can more effectively explore long-range dependencies. Moreover, treating relational structure as a constraint on inference paths naturally supports retrieval governed by exact rules and symbolic constraints. Representative approaches such as AriGraph (Anokhin et al., 2024), EMG-RAG (Wang et al., 2024k), Mem0$^g$ (Chhikara et al., 2025), and SGMem (Wu et al., 2025h) first identify the most relevant nodes or triples and then expand to their semantically related K-hop neighbors to construct an ego-graph. HippoRAG (Gutierrez et al., 2024) performs personalized PageRank (Page et al., 1999) seeded on the retrieved nodes and ranks the rest of the graph by their proximity to these seeds, enabling effective multi-hop retrieval. Going beyond fixed expansion rules, CAM (Li et al., 2025f) and D-SMART (Lei et al., 2025) employ LLMs to steer subgraph exploration: CAM uses an LLM to select informative neighbors and children of a central node for associative exploration, while D-SMART treats the LLM as a planner that performs beam search over a KG memory to retrieve one-hop neighbors of target entities and the relations connecting a given entity pair. For temporal graphs, Zep (Rasmussen et al., 2025) and MemoTime (Tan et al., 2025b) further enable entity-subgraph construction and relation retrieval under explicit temporal constraints, ensuring that the returned results satisfy the required time rules.

**Generative Retrieval**   This strategy replaces lexical or semantic retrieval with a model that directly generates the identifiers of relevant documents (Tay et al., 2022; Wang et al., 2022b). By framing retrieval as a conditional generation task, the model implicitly stores candidate documents in its parameters and performs deep query–document interaction during decoding(Li et al., 2025j). Leveraging the semantic capabilities of pretrained language models, this paradigm often outperforms traditional retrieval methods, particularly in small-scale settings(Zeng et al., 2024). However, generative retrieval requires additional training to internalize the semantics of all candidate documents, resulting in limited scalability when the corpus evolves (Yuan et al., 2024b). For these reasons, agentic memory systems have paid relatively little attention to this paradigm, although its tight integration of generation and retrieval suggests untapped potential.

**Hybrid Retrieval**   This strategy integrates the strengths of multiple retrieval paradigms. Systems such as Agent KB (Tang et al., 2025d) and MIRIX (Wang and Chen, 2025) combine lexical and semantic retrieval to balance precise term or tool matching with broader semantic alignment. Similarly, Semantic Anchoring (Chatterjee and Agarwal, 2025) performs parallel searches over semantic embeddings and symbolic inverted indices to achieve complementary coverage. Some other methods combine multiple evaluation signals to guide retrieval. Generative Agents (Kaiya et al., 2023), for example, illustrate this multi-factor approach through a scoring scheme that accumulates recency, importance, and relevance. MAICC (Jiang et al., 2025c) adopts a mixed-utility scoring function that integrates similarity with both global and predicted individual returns. In graph-based settings, retrieval typically proceeds in two stages: semantic retrieval first identifies relevant nodes or triples, and graph topology is subsequently leveraged to expand the search space (Anokhin et al., 2024; Wang et al., 2024k; Gutierrez et al., 2024; Li et al., 2025f).

At the database infrastructure level, MemoriesDB (Ward, 2025) introduces a temporal–semantic–relational database designed for long-term agent memory, providing a hybrid retrieval architecture that integrates these dimensions into a unified storage and access framework.

By fusing heterogeneous retrieval signals, hybrid approaches preserve the precision of keyword matching while incorporating the contextual understanding of semantic methods, ultimately yielding more comprehensive and relevant results.

### 5.3.4   Post-Retrieval Processing

Initial retrieval often returns fragments that are redundant, noisy, or semantically inconsistent. Directly injecting these results into the prompt can lead to excessively long contexts, conflicting information, and reasoning distracted by irrelevant content. Post-retrieval processing, therefore, becomes essential for ensuring prompt quality. Its goal is to distill the retrieved results into a concise, accurate, and semantically coherent context. In practice, two components are central: (1) **Re-ranking and Filtering:** performing fine-grained relevance estimation to remove irrelevant or outdated memories and reorder the remaining fragments, thereby reducing noise and redundancy. (2) **Aggregation and Compression:** integrating the retrieved memories with the original query, eliminating duplication, merging semantically similar information, and reconstructing a compact and coherent final context.

**Re-ranking and Filtering**   To maintain a concise and coherent context, initial retrieval results are re-ranked and filtered to ensure a concise and coherent context by removing low-relevance items. Early approaches rely on heuristic criteria for evaluating semantic consistency. For example, Semantic Anchoring (Chatterjee and Agarwal, 2025) integrates vector similarity with entity- and discourse-level alignment, whereas RCR-Router (Liu et al., 2025c) combines multiple handcrafted signals, including role relevance, task-stage priority, and recency. These methods, however, often require extensive hyperparameter tuning to balance heterogeneous importance scores. To alleviate this burden, learn-to-memorize (Zhang et al., 2025t) formulates score aggregation as a reinforcement-learning problem, enabling the model to learn optimal weights over retrieval signals. While these techniques primarily optimize semantic coherence, scenarios demanding strict temporal reasoning require additional constraints: Rasmussen et al. (2025) and Tan et al. (2025b) filter memories based on their timestamps and validity windows to satisfy complex temporal dependencies.

With the increasing capability of LLMs, recent methods leverage their intrinsic language understanding to assess memory quality directly. Memory-R1 (Yan et al., 2025b) and Westhäußer et al. (2025) both introduce

LLM-based evaluators (Answer Agents or Self-Validator Agents) that filter retrieved content before producing the final response. However, prompt-based filtering remains limited by the LLM's inherent capacity and by mismatches between prompt semantics and downstream usage. Consequently, many systems train auxiliary models to estimate memory importance more robustly (Tan et al., 2025c). Memento (Zhou et al., 2025a) uses Q-learning (Watkins and Dayan, 1992) to predict the probability that a retrieved item contributes to a correct answer, and MemGuide (Du et al., 2025b) fine-tunes LLaMA-8B (Grattafiori et al., 2024) to re-rank candidates using marginal slot-completion gain. Together, these re-ranking and filtering strategies refine retrieval results without modifying the underlying retriever, enabling compatibility with any pre-trained retrieval model while supporting task-specific optimization.

**Aggregation and Compression**  Another approach to improving both the quality and efficiency of downstream reasoning through post-retrieval processing is the aggregation and compression. This process integrates the retrieved evidence with the query to form a coherent and compact context. Unlike filtering and re-ranking, which mainly address noise and prioritization, this stage focuses on merging multiple fragmented memory items into higher-level and distilled knowledge representations, and on refining these representations when task-specific adaptations are required. ComoRAG (Wang et al., 2025f) illustrates this idea through its Integration Agent, which identifies historical signals that are semantically aligned with the query and combines them into an abstract global summary that provides broad contextual grounding. The Extractor Agent in MA-RAG (Nguyen et al., 2025) performs fine-grained content selection over the retrieved documents, retaining only the key information that is strongly relevant to the current subquery and producing concise snippets tailored to local reasoning needs.

Furthermore, G-Memory (Zhang et al., 2025c) extends aggregation and compression into the personalization for multi-agent systems. It consolidates retrieved high-level insights and sparsified trajectories, and then uses an LLM to customize these condensed experiences according to the agent's role. This process refines general knowledge into role-specific prompts that populate the agent's personalized memory.

**Summary**  In conclusion, post-retrieval processing acts as a crucial intermediate step that transforms noisy, fragmented retrieval results into a precise and coherent context for reasoning. Through above mechanisms, the post-retrieval processing not only enhances the density and fidelity of the memories supplied to the model but also aligns the information with task requirements and agent characteristics.

# 6  Resources and Frameworks

## 6.1  Benchmarks and Datasets

In this section, we survey representative benchmarks and datasets that have been used (or could be used) to evaluate the memory, long-term, continual-learning, or long-context capabilities of LLM-based agents. We classify these benchmarks into two broad categories: (1) those explicitly designed for memory / lifelong learning / self-evolving agents, and (2) those originally developed for other purposes (e.g., tool-use capacity, web search, embodied action) but nevertheless relevant for memory evaluation due to their long-horizon, multi-task, or sequential nature.

### 6.1.1  Benchmarks for Memory / Lifelong / Self-Evolving Agents

Memory-oriented benchmarks focus primarily on how well an agent can construct, maintain, and exploit an explicit memory of past interactions or world facts. These tasks typically probe the retention and retrieval of information across multi-turn dialogues, user-specific sessions, or long synthetic narratives, sometimes including multimodal signals.

A consolidated overview of these benchmarks, including their memory focus, environment type, modality, and evaluation scale, is provided in Table 8, which serves as a structured reference for comparing their design objectives and evaluation settings. Representative examples such as MemBench (Tan et al., 2025a), LoCoMo (Maharana et al., 2024), WebChoreArena (Miyai et al., 2025), MT-Mind2Web (Deng et al., 2024), PersonaMem (Jiang et al., 2025a), PerLTQA (Du et al., 2024), MPR (Zhang et al., 2025u), PrefEval (Zhao

**Table 8** Overview of benchmarks relevant to LLM agent memory, long-term, lifelong learning, and self-evolving evaluation. The table covers two categories of benchmarks: (i) benchmarks explicitly designed for memory-, lifelong learning-, or self-evolving agent evaluation, and (ii) other agent-oriented benchmarks that implicitly stress long-horizon memory through sequential, multi-step, or multi-task interactions. **Fac.** and **Exp.** indicate whether a benchmark evaluates factual memory or experiential (interaction-derived) memory, respectively. **MM.** denotes the presence of multimodal inputs, while **Env.** indicates whether the benchmark is conducted in a simulated or real environment. **Feature** summarizes the primary capability under evaluation, and **Scale** reports the approximate benchmark size in terms of *samples* (s.) or *tasks* (t.). PDDL denotes commonly used PDDL-based planning subsets.

| Name | Link | Fac. | Exp. | MM. | Env. | Feature | Scale |
|---|---|---|---|---|---|---|---|
| **Memory/Lifelong-learning/Self-evolving-oriented Benchmarks** | | | | | | | |
| MemBench | GitHub | ✔ | ✔ | ✗ | simulated | interactive scenarios | 53,000 s. |
| MemoryAgentBench | GitHub | ✔ | ✔ | ✗ | simulated | multi-turn interactions | 4 t. |
| LoCoMo | Website | ✔ | ✗ | ✔ | real | conversational memory | 300 s. |
| WebChoreArena | GitHub | ✔ | ✔ | ✔ | real | tedious web browsing | 4 t./532 s. |
| MT-Mind2Web | GitHub | ✔ | ✔ | ✗ | real | conversational web navigation | 720 s. |
| PersonaMem | Website | ✔ | ✗ | ✗ | simulated | dynamic user profiling | 15 t./180 s. |
| LongMemEval | GitHub | ✔ | ✗ | ✗ | simulated | interactive memory | 5 t./500 s. |
| PerLTQA | Website | ✔ | ✗ | ✗ | simulated | social personalized interactions | 8,593 s. |
| MemoryBank | Website | ✔ | ✗ | ✗ | simulated | user memory updating | 194 s. |
| MPR | GitHub | ✔ | ✗ | ✗ | simulated | user personalization | 108,000 s. |
| PrefEval | Website | ✔ | ✗ | ✗ | simulated | personal preferences | 3,000 s. |
| LOCCO | Website | ✔ | ✗ | ✗ | simulated | chronological conversations | 3,080 s. |
| StoryBench | Website | ✔ | ✔ | ✗ | mixed | interactive fiction games | 3 t. |
| MemoryBench | Website | ✔ | ✔ | ✗ | simulated | continual learning | 4 t./∼ 20,000 s. |
| Madial-Bench | GitHub | ✔ | ✗ | ✗ | simulated | memory recalling | 331 s. |
| Evo-Memory | Website | ✔ | ✔ | ✗ | simulated | test-time learning | 10 t./∼ 3,700 s. |
| LifelongAgentBench | Website | ✔ | ✔ | ✗ | simulated | lifelong learning | 1,396 s. |
| StreamBench | Website | ✔ | ✔ | ✗ | simulated | continuous online learning | 9,702 s. |
| DialSim | Website | ✔ | ✔ | ✗ | real | multi-dialogue understanding | ∼ 1,300 s. |
| LongBench | Website | ✔ | ✗ | ✗ | mixed | long-context understanding | 21 t./4,750 s. |
| LongBench v2 | Website | ✔ | ✗ | ✗ | mixed | long-context multitasks | 20 t./503 s. |
| RULER | GitHub | ✔ | ✗ | ✗ | simulated | long-context retrieval | 13 t. |
| BABILong | GitHub | ✔ | ✗ | ✗ | simulated | long-context reasoning | 20 t. |
| MM-Needle | Website | ✔ | ✗ | ✔ | simulated | multimodal long-context retrieval | ∼ 280,000 s. |
| HaluMem | GitHub | ✔ | ✗ | ✗ | simulated | memory hallucinations | 3,467 s. |
| HotpotQA | Website | ✔ | ✗ | ✗ | simulated | long-context QA | 113k s. |
| **Other Related Benchmarks** | | | | | | | |
| ALFWorld | Website | ✔ | ✔ | ✗ | simulated | text-based embodied environment | 3,353 t. |
| ScienceWorld | GitHub | ✔ | ✔ | ✗ | simulated | interactive embodied environment | 10 t./30 t. |
| AgentGym | Website | ✗ | ✔ | ✗ | mixed | multiple environments | 89 t./20,509 s. |
| AgentBoard | GitHub | ✗ | ✔ | ✗ | mixed | multi-round interaction | 9 t./1013 s. |
| PDDL* | Website | ✗ | ✔ | ✗ | simulated | strategy game | - |
| BabyAI | Website | ✗ | ✔ | ✗ | simulated | language learning | 19 t. |
| WebShop | Website | ✗ | ✔ | ✔ | simulated | e-commerce web interaction | 12,087 s. |
| WebArena | Website | ✗ | ✔ | ✔ | real | web interaction | 812 s. |
| MMInA | Website | ✔ | ✔ | ✔ | real | multihop web interaction | 1,050 s. |
| SWE-Bench Verified | Website | ✗ | ✔ | ✗ | real | code repair | 500 s. |
| GAIA | Website | ✗ | ✔ | ✔ | real | human-level deep research | 466 s. |
| xBench-DS | Website | ✗ | ✔ | ✔ | real | deep-search evaluation | 100 s. |
| ToolBench | GitHub | ✗ | ✔ | ✗ | real | API tool use | 126,486 s. |
| GenAI-Bench | Website | ✗ | ✔ | ✔ | real | visual generation evaluation | ∼ 40,000 s. |

et al., 2025d), LOCCO (Jia et al., 2025), StoryBench (Wan and Ma, 2025), Madial-Bench (He et al., 2025),

DialSim (Zheng et al., 2025b), LongBench (Bai et al., 2024), LongBench v2 (Bai et al., 2025), RULER (Hsieh et al., 2024), BALILong (Kuratov et al., 2024) MM-Needle (Wang et al., 2025e), and HaluMem (Packer et al., 2023a) stress user modeling, preference tracking, and conversation-level consistency, often under simulated settings where ground-truth memories can be precisely controlled.

Lifelong-learning benchmarks extend beyond isolated memory retrieval to examine how agents continually acquire, consolidate, and update knowledge over long horizons and evolving task distributions. Benchmarks such as LongMemEval (Wu et al., 2025a), MemoryBank (Zhong et al., 2024), MemoryBench (Ai et al., 2025), LifelongAgentBench (Zheng et al., 2025b), and StreamBench (Wu et al., 2024a), are designed around sequences of tasks or episodes in which new information gradually arrives and earlier information may become obsolete or conflicting. These setups emphasize phenomena like catastrophic forgetting, forward and backward transfer, and test-time adaptation, making them suitable for studying how memory mechanisms interact with continual-learning objectives. In many cases, performance is tracked not only on the current task but also on previously seen tasks or conversations, thereby quantifying how well the agent preserves useful knowledge while adapting to new users, domains, or interaction patterns.

Self-evolving-agent benchmarks go a step further by treating the agent as an open-ended system that can iteratively refine its own memory, skills, and strategies through interaction. Here, the focus is not only on storing and recalling information, but also on meta-level behaviors such as self-reflection, memory editing, tool-augmented storage, and policy improvement over multiple episodes or games. Benchmarks like MemoryAgentBench (Hu et al., 2025c), Evo-Memory (Wei et al., 2025e), and other multi-episode or mission-style environments can be instantiated in a self-evolving setting by allowing the agent to accumulate trajectories, synthesize higher-level abstractions, and adjust its behavior in future runs based on its own past performance. When viewed through this lens, these benchmarks provide a testbed for evaluating whether an agent can autonomously bootstrap more capable behaviors over time-turning static tasks into arenas for long-term adaptation, strategy refinement, and genuinely self-improving memory use.

### 6.1.2 Other Related Benchmarks

Beyond benchmarks explicitly designed for memory or lifelong learning, a wide range of agent-oriented and long-horizon evaluation suites are also relevant for studying memory-related capabilities in LLM-based agents. Although these benchmarks were originally introduced to assess other aspects such as tool use, embodied interaction, or knowledge-intensive reasoning, their sequential, multi-step, and multi-task nature implicitly places strong demands on long-term information retention, context management, and state tracking.

Embodied and interactive environments constitute a major class of such benchmarks. Frameworks like ALFWorld (Shridhar et al., 2021) and ScienceWorld (Wang et al., 2022a) evaluate agents in simulated text-based or partially grounded environments where success requires remembering past observations, intermediate goals, and environment dynamics across extended action sequences. Similarly, BabyAI (Chevalier-Boisvert et al., 2019) focuses on language-conditioned instruction following over temporally extended episodes, implicitly testing an agent's ability to maintain task-relevant state throughout interaction. While these benchmarks do not explicitly model external memory modules, effective performance often depends on the agent's capacity to preserve and reuse information over long horizons.

Another prominent category includes web-based and tool-augmented interaction benchmarks. WebShop (Yao et al., 2023a), WebArena (Zhou et al., 2024b), and MMInA (Tian et al., 2025) assess agents operating in realistic or semi-realistic web environments involving multi-step navigation, information gathering, and decision making. These settings naturally induce long-context trajectories in which earlier actions, retrieved information, or user constraints must be recalled and integrated at later stages. ToolBench (Qin et al., 2024a) further extends this paradigm by evaluating an agent's ability to select and invoke APIs across complex workflows, where memory of prior tool outputs and tool-use experience is critical for coherent execution.

Multi-task and general agent evaluation platforms also provide indirect but valuable signals about memory usage. AgentGym (Xi et al., 2024b) and AgentBoard (Xi et al., 2024b) aggregate diverse environments or tasks into unified evaluation suites, requiring agents to adapt across tasks while retaining task-specific knowledge and strategies. PDDL-based planning environments, commonly used in agent benchmarks, evaluate strategic reasoning over structured action spaces, where agents benefit from accumulating and reusing experience across episodes to improve long-horizon planning performance.

**Table 9** Overview of representative open-source memory frameworks for LLM-based agents. The table compares widely used frameworks in terms of the types of memory they support (factual vs. experiential), multimodality, internal memory structure, and reported evaluation benchmarks. **Fac.** and **Exp.** denote factual and experiential memory, respectively, **MM.** indicates multimodal memory support, and **Structure** summarizes the core memory abstraction or organization mechanism adopted by each framework. **Evaluation** lists publicly reported benchmarks used to assess memory-related capabilities, when available.

| Framework | Links | Fac. | Exp. | MM. | Structure | Evaluation |
|---|---|---|---|---|---|---|
| MemGPT | GitHub Website | ✔ | ✔ | ✘ | hierachical (S/LTM) | LoCoMo |
| Mem0 | GitHub Website | ✔ | ✔ | ✘ | graph + vector | LoCoMo |
| Memobase | GitHub Website | ✔ | ✔ | ✘ | structured profiles | LoCoMo |
| MIRIX | GitHub Website | ✔ | ✔ | ✔ | structured memory | LoCoMo, MemoryAgentBench |
| MemoryOS | GitHub Website | ✔ | ✔ | ✘ | hierarchical (S/M/LTM) | LoCoMo, MemoryBank |
| MemOS | GitHub Website | ✔ | ✔ | ✘ | tree memory + memcube | LoCoMo, PreFEval, LongMemEval, PersonaMem |
| Zep | GitHub Website | ✔ | ✔ | ✘ | temporal knowledge graph | LongMemEval |
| LangMem | GitHub Website | ✔ | ✔ | ✘ | core API + manager | - |
| SuperMemory | GitHub Website | ✔ | ✔ | ✔ | vector + semantic | - |
| Cognee | GitHub Website | ✔ | ✔ | ✔ | knowledge graph | - |
| Memary | GitHub Website | ✔ | ✔ | ✘ | stream + entity store | - |
| Pinecone | GitHub Website | ✔ | ✘ | ✘ | vector database | - |
| Chroma | GitHub Website | ✔ | ✘ | ✔ | vector database | - |
| Weaviate | GitHub Website | ✔ | ✘ | ✔ | vector + graph | - |
| Second Me | GitHub Website | ✔ | ✘ | ✘ | agent ego | - |
| MemU | GitHub Website | ✔ | ✔ | ✔ | hierachical layers | - |
| MemEngine | GitHub | ✔ | ✔ | ✔ | modular space | - |
| Memori | GitHub Website | ✔ | ✔ | ✘ | memory database | - |
| ReMe | GitHub Website | ✔ | ✔ | ✘ | memory management | - |
| AgentMemory | GitHub Website | ✔ | ✔ | ✘ | memory management | - |
| MineContext | GitHub Website | ✔ | ✔ | ✔ | context engineering | - |
| Acontext | GitHub | ✔ | ✔ | ✔ | context engineering + skill learning | - |

Finally, several recent benchmarks target demanding real-world or near-real-world reasoning scenarios that inherently stress long-context and cross-step consistency. SWE-Bench Verified (Jimenez et al., 2024) evaluates code repair over realistic software repositories, where agents must reason over long files and evolving code states. GAIA (Mialon et al., 2023) and xBench (Chen et al., 2025b) assess deep research and search-intensive tasks that require synthesizing information gathered across multiple steps and sources. GenAI-Bench (Li et al., 2024a), while focusing on multimodal generation quality, similarly involves complex workflows in which memory of prior prompts, intermediate outputs, or visual constraints plays a nontrivial role.

Taken together, these benchmarks complement memory-oriented evaluations explicitly by situating LLM-based agents in rich, interactive, and long-horizon settings. Although memory is not always an explicit target of measurement, sustained performance in these environments implicitly depends on an agent's ability to manage long contexts, preserve relevant information, and integrate past experience into ongoing decision making, making them valuable testbeds for studying memory-related behaviors in practice.

## 6.2 Open-Source Frameworks

A rapidly growing ecosystem of open-source memory frameworks aims to provide reusable infrastructure for building memory-augmented LLM agents. A structured comparison of representative open-source memory frameworks, including their supported memory types, architectural abstractions, and evaluation coverage, is

summarized in Table 9. Most of these frameworks support factual memory via vector or structured stores, and an increasing subset also models experiential traces, such as dialogue histories, user actions, and episodic summaries, with multimodal memory emerging more recently. Open-source memory frameworks for LLM agents span a spectrum from agent-centric systems with rich, hierarchical memory abstractions to more general-purpose retrieval or memory-as-a-service backends, e.g., MemGPT (Packer et al., 2023b), Mem0 (Chhikara et al., 2025), Memobase, MemoryOS (Kang et al., 2025a), MemOS (Li et al., 2025k), Zep (Rasmussen et al., 2025), LangMem (LangChain, 2025), SuperMemory (Supermemory, 2025), Cognee (Cognee, 2025), Memary (Memary, 2025), Pinecone, Chroma, Weaviate, Second Me, MemU, MemEngine (Zhang et al., 2025s), Memori, ReMe (AgentScope, 2025), AgentMemory, and MineContext (MineContext, 2025). Many of them explicitly separate short- and long-term stores and offer graph-based, profile-based, or modular memory spaces, and some have begun to report results on memory-based benchmarks. The others typically provide scalable vector or graph databases, APIs, and semantic or streaming entity layers that help organize context but often leave agent behavior and evaluation protocols to the application. Overall, these frameworks are rapidly maturing in their representational flexibility and system design.

# 7 Positions and Frontiers

This section articulates key positions and emerging frontiers in the design of memory systems for LLM-based agents. Moving beyond descriptive surveys of existing methods, we focus on paradigm-level shifts that redefine how memory is constructed, managed, and optimized in long-horizon agentic settings. Specifically, we examine the transition from retrieval-centric to generative memory, from manually engineered to autonomously managed memory systems, and from heuristic pipelines to reinforcement learning–driven memory control. We further discuss how these shifts intersect with multimodal reasoning, multi-agent collaboration, and trustworthiness, outlining open challenges and research directions that are likely to shape the next generation of agent memory architectures.

## 7.1 Memory Retrieval vs. Memory Generation

### 7.1.1 Look Back: From Memory Retrieval to Memory Generation

Historically, the dominant paradigm in agent memory research has centered on **memory retrieval**. Under this paradigm, the primary objective is to identify, filter, and select the most relevant memory entries from an existing memory store given the current context. A large body of prior work focuses on improving retrieval accuracy through better indexing strategies, similarity metrics, reranking models, or structured representations such as knowledge graphs (Tan et al., 2025c; Memobase, 2025). In practice, this includes techniques such as vector similarity search with dense embeddings, hybrid retrieval combining lexical and semantic signals, hierarchical filtering, and graph-based traversal. These methods emphasize precision and recall in accessing stored information, implicitly assuming that the memory base itself is already well formed.

Recently, however, increasing attention has shifted toward **memory generation**. Rather than treating memory as a static repository to be queried, memory generation emphasizes the agent's ability to actively synthesize new memory representations on demand. The goal is not merely to retrieve and concatenate existing fragments, but to integrate, compress, and reorganize information in a manner that is tailored to the current context and future utility. This shift reflects a growing recognition that effective memory usage often requires abstraction and recomposition, especially when raw stored information is noisy, redundant, or misaligned with the immediate task.

Existing approaches to memory generation can be broadly grouped into two directions. One line of work adopts a **retrieve then generate** strategy, where retrieved memory items serve as raw material for reconstruction. In this setting, the agent first accesses a subset of relevant memories and then generates a refined memory representation that is more concise, coherent, and context specific, as implemented in ComoRAG (Wang et al., 2025f), G-Memory (Zhang et al., 2025c) and CoMEM (Wu et al., 2025d). This approach preserves grounding in historical information while enabling adaptive summarization and restructuring. A second line of work explores **direct memory generation**, in which memory is produced without any explicit retrieval step. Instead, the agent generates memory representations directly from the current context, interaction history, or latent internal states. Systems such as MemGen (Zhang et al., 2025d) and VisMem (Yu et al., 2025e) exemplify this

direction by constructing latent memory tokens that are customized to the task at hand, bypassing explicit memory lookup altogether.

### 7.1.2 Future Perspective

Looking ahead, we anticipate that generative approaches will play an increasingly central role in agent memory systems. We highlight three properties that future generative memory mechanisms should ideally exhibit.

First, generative memory should be **context adaptive**. Rather than storing generic summaries, the memory system should generate representations that are explicitly optimized for the agent's anticipated future needs. This includes adapting the granularity, abstraction level, and semantic focus of memory to different tasks, stages of problem solving, or interaction regimes.

Second, generative memory should support **integration across heterogeneous signals**. Agents increasingly operate over diverse modalities and information sources, including text, code, tool outputs, and environmental feedback. Memory generation provides a natural mechanism for fusing these fragmented signals into unified representations that are more useful for downstream reasoning than raw concatenation or retrieval alone. We hypothesize that latent memory (as discussed in Section 3.3) might be a promising technical path for this gaol.

Third, generative memory should be **learned and self optimizing**. Rather than relying on manually specified generation rules, future systems should learn when and how to generate memory through optimization signals, such as reinforcement learning or long horizon task performance. In this view, memory generation becomes an integral component of the agent's policy, co evolving with reasoning and decision making.

## 7.2 Automated Memory Management

### 7.2.1 Look-Back: From Hand-crafted to Automatically Constructed Memory Systems.

Existing agent memory systems (Xu et al., 2025c; Packer et al., 2023a) typically rely on manually designed strategies to determine what information to store, when to use it, and how to update or retrieve it. By guiding fixed LLMs with detailed instructions (Chhikara et al., 2025), predefined thresholds (Kang et al., 2025a), or explicit human-crafted rules drafted by human experts (Xu et al., 2025c), system designers can integrate memory modules into current agent frameworks with relatively low computational and engineering cost, enabling rapid prototyping and deployment. Besides, they also offer **interpretability, reproducibility, and controlled**, allowing the developers to precisely specify the state and behavior of memory. However, similar to expert systems in other areas, such manually curated approaches suffer from significant limitations: they are inherently inflexible and often fail to generalize across diverse, dynamic environments. Consequently, these systems tend to underperform in long-term or open-ended interactions.

Recent developments in agent memory research begin to address these limitations by enabling the agents themselves to autonomously manage the memory evolution and retrieval. For example, CAM (Li et al., 2025f) empowers LLM agents to automatically cluster fine-grained memory entries into high-level abstract units. Memory-R1 (Yan et al., 2025b) introduces an auxiliary agent equipped with a dedicated "memory manager" tool to handle memory updates. Despite these advances, current solutions remain constrained: many are still driven by manually engineered rules or are optimized for narrow, task-specific learning objectives, making them difficult to generalize to open-ended settings.

### 7.2.2 Future Perspective

To support truly automated memory management, a promising direction is to **integrate memory construction, evolution, and retrieval directly into the agent's decision loop via explicit tool calls**, making the agent itself reason about memory operations instead of depending on external modules or hand-crafted workflows. Compared with existing designs that separate an agent's internal reasoning process from its memory management actions, an LLM agent can know precisely what memory actions it performs (e.g., add/update/delete/retrieval) in this tool-based strategy, leading to more coherent, transparent, and contextually grounded memory behavior.

Another key frontier lies in developing **self-optimizing memory structures** adopting hierarchical and adaptive architectures inspired by cognitive systems. First, hierarchical memory structure has been shown to improve

the efficiency and performance (Kang et al., 2025a). Beyond hierarchy, self-evolving memory systems that dynamically link, index, and reconstruct memory entries enable the memory storage itself to self-organize over time, supporting richer reasoning and reducing dependence on hand-designed rules. Ultimately, such adaptive, self-organizing memory architectures pave the way toward agents capable of maintaining robust, scalable, and truly autonomous memory management.

## 7.3 Reinforcement Learning Meets Agent Memory



**Figure 11** The evolution of RL-enabled agent memory systems. A conceptual progression from **RL-free** memory systems based on heuristic or prompt-driven pipelines, to **partially RL-involved** designs where reinforcement learning governs selected memory operations, and finally to fully **RL-driven** memory systems in which memory architectures and control policies are learned end-to-end. This evolution reflects a broader paradigm shift from *manually* engineered memory pipelines toward *model-native*, *self-optimizing* memory management in LLM-based agents.

### 7.3.1 Look-Back: RL is Internalizing Memory Management Abilities for Agents.

Reinforcement learning is rapidly reshaping the development paradigm of modern LLM-based agents. Across a wide spectrum of agentic capabilities, including planning, reasoning, tool use, as well as across diverse task domains such as mathematical reasoning, deep research, and software engineering, RL has begun to play a central role in driving agent performance (Zhang et al., 2025f,k). Memory, as one of the foundational components of agentic capability, follows a similar trend from pipeline-based to model-native paradigm (Sang et al., 2025). The agent memory research community is collectively transitioning from **early heuristic and manually engineered designs** to approaches in which **RL increasingly governs key decisions**. Looking ahead, it is reasonable to expect that **fully RL–based memory systems** may eventually become the dominant direction. Before discussing this trajectory in detail, we briefly outline the first stage of development. This transition, in which memory management is progressively internalized and optimized through reinforcement learning, is schematically illustrated in Figure 11.

**RL-free Memory Systems** A substantial portion of the agent memory literature surveyed earlier can be categorized as RL-free memory systems. These approaches typically rely on heuristic or manually specified mechanisms, such as fixed thresholding rules inspired by curves of forgetting, rigid semantic search pipelines found in frameworks such as MemOS (Li et al., 2025k), Mem0 (Chhikara et al., 2025), and MemoBase (Memobase, 2025), or simple concatenation-based strategies for storing memory chunks. In some systems, an LLM participates in memory management in a way that appears *agentic*, yet the underlying behavior is entirely prompt-driven. The LLM is asked to generate memory entries but has not received any dedicated training for effective memory control, as seen in systems such as Dynamic Cheatsheet (Suzgun et al., 2025), ExpeL (Zhao et al., 2024), EvolveR (Wu et al., 2025c), and G-Memory (Zhang et al., 2025c). This class of methods has dominated early work in the field and is likely to remain influential for some time due to its simplicity and practical accessibility.

**RL-assisted Memory Systems** As the field progressed, many works began to incorporate RL-based methods into selected components of the memory pipeline. An early attempt in this direction is RMM (Tan et al., 2025c), which employed a lightweight policy gradient learner to rank memory chunks after an initial retrieval stage based on BM25 or other semantic similarity metrics. Later systems explored substantially more ambitious designs. For example, Mem-$\alpha$ (Wang et al., 2025o) delegates the entire process of memory construction to an agent trained with RL, and Memory-R1 (Yan et al., 2025b) employs a similar philosophy. A rapidly expanding line of research investigates how an agent can autonomously fold, compress, and manage context in ultra-long multi-turn tasks. This setting corresponds to the management of working memory (Kang et al., 2025c; Ye et al., 2025a). Many of the leading systems in this area are trained with RL, including but not limited to Context Folding (Sun et al., 2025a), Memory-as-Action (Zhang et al., 2025q), MemSearcher (Yuan et al., 2025a), and IterResearch (Chen et al., 2025a). These RL-assisted approaches have already demonstrated strong capabilities and point toward the increasing role of RL in future memory system design.

### 7.3.2 Future Perspective

Looking forward, we anticipate that **fully RL-driven memory systems will constitute the next major stage** in the evolution of agent memory. We highlight two properties that such systems should ideally embody.

- First, *memory architectures managed by agents should minimize reliance on human-engineered priors.* Many existing frameworks inherit design patterns inspired by human cognition, such as cortical or hippocampal analogies (Gutierrez et al., 2024), or predefined hierarchical taxonomies that partition memory into episodic, semantic, and core categories (Wang and Chen, 2025). Although these abstractions have been useful for grounding early work, they may not represent the most effective or natural structures for artificial agents operating in complex environments. A fully RL-driven setting offers the possibility for agents to invent novel and potentially more suitable memory organizations that emerge directly from optimization dynamics rather than human intuition. In this view, the agent is encouraged to design new memory formats, storage schemas, or update rules through RL incentives, enabling memory architectures that are adaptive and creative rather than handcrafted.

- Second, *future memory systems should provide agents with complete control over all stages of memory management.* Current RL-assisted approaches typically intervene in only a subset of the memory lifecycle. For instance, Mem-$\alpha$ automates certain aspects of memory writing yet still relies on manually defined retrieval pipelines, whereas systems such as MemSearcher (Yuan et al., 2025a) focus primarily on short-term working memory without addressing long-term consolidation or evolution. A fully agentic memory system would require the agent to autonomously handle multi-granular memory formation, memory evolution, and memory retrieval in an integrated manner. Achieving this level of control will almost certainly require end-to-end RL training, since heuristic or prompt-based methods are insufficient for coordinating the complex interactions among these components across long-time horizons.

Together, these two directions suggest a future in which memory is not merely an auxiliary mechanism bolted onto an LLM agent, but rather a fully learnable and self-organizing subsystem that coevolves with the agent through RL. Such systems hold the potential to enable genuinely continual learning and long-term competence in artificial agents.

## 7.4 Multimodal Memory

### 7.4.1 Look-Back

As research on text-based memory becomes increasingly mature and extensively explored, and as multimodal large language models and unified models that jointly support multimodal understanding and generation continue to advance, attention has naturally expanded toward **multimodal memory**. This shift reflects a broader recognition that many real-world agentic settings are inherently multimodal, and that memory systems limited to text alone are insufficient to support long-horizon reasoning and interaction in complex environments.

Existing efforts on multimodal memory can be broadly grouped into two complementary directions. The first focuses on enabling **multimodal agents** to store, retrieve, and utilize memories derived from diverse sensory inputs (Long et al., 2025; Zuo et al., 2025). This direction is a natural extension of agent memory, since agents

operating in realistic environments inevitably encounter heterogeneous data sources, including images, audio, video, and other non-textual signals (Xie et al., 2024). The degree of progress in multimodal memory closely follows the maturity of corresponding modalities. Visual modalities such as images and videos have received the most attention, leading to a growing body of work on visual and video memory mechanisms that support tasks such as visual grounding, temporal tracking, and long-term scene consistency (Long et al., 2025; Wang et al., 2024g; Gurukar and Kadav, 2025; Yu et al., 2025e; Bo et al., 2025; Wang et al., 2025p; Li et al., 2024d). In contrast, memory systems for audio and other modalities remain relatively underexplored (Li et al., 2025a).

The second direction treats memory as an enabling component for **unified models**. In this setting, memory is leveraged not primarily to support agent decision making, but to enhance multimodal generation and consistency. For example, in image and video generation systems, memory mechanisms are often used to preserve entity consistency, maintain world state across frames, or ensure coherence across long generation horizons (Yu et al., 2025b). Here, memory serves as a stabilizing structure that anchors generation to previously produced content, rather than as a record of agent experience per se.

### 7.4.2 Future Perspective

Looking forward, multimodal memory is likely to become an indispensable component of agentic systems. As agents increasingly move toward embodied and interactive settings, their information sources will be inherently multimodal, spanning perception, action, and environmental feedback. Effective memory systems must therefore support the storage, integration, and retrieval of heterogeneous signals in a unified manner.

Despite recent progress, there is currently no memory system that provides truly **omnimodal support**. Most existing approaches remain specialized to individual modalities or loosely coupled across modalities. A key future challenge lies in designing memory representations and operations that can flexibly accommodate diverse modalities while preserving semantic alignment and temporal coherence. Moreover, multimodal memory must evolve beyond passive storage to support abstraction, cross-modal reasoning, and long-term adaptation. Addressing these challenges will be essential for enabling agents that can operate robustly and coherently in rich, multimodal environments.

## 7.5 Shared Memory in Multi-Agent Systems

### 7.5.1 Look-Back: From Isolated Memories to Shared Cognitive Substrates

As LLM-based multi-agent systems (MAS) have gained prominence, **shared memory** has emerged as a key mechanism for enabling coordination, consistency, and collective intelligence. Early multi-agent frameworks primarily relied on **isolated local memories** coupled with explicit message passing, where agents exchanged information through dialogue histories or task-specific communication protocols (Qian et al., 2024; Wu et al., 2024b; Hu et al., 2025b; Zhang et al., 2025i). While this design avoided direct interference between agents, it often suffered from redundancy, fragmented context, and high communication overhead, especially as team size and task horizon increased.

Subsequent work introduced **centralized shared memory structures**, such as global vector stores, blackboard systems, or shared documents (Hong et al., 2024), accessible to all agents. These designs enabled a form of team-level memory that supported joint attention, reduced duplication, and facilitated long-horizon coordination. Representative systems demonstrated that shared memory could serve as a persistent common ground for planning, role handoff, and consensus building (Rezazadeh et al., 2025b; Xu et al., 2025a). However, naive global sharing also exposed new challenges, including memory clutter, write contention, and the lack of role- or permission-aware access control.

### 7.5.2 Future Perspective

Looking forward, shared memory is likely to evolve from a passive repository into an **actively managed and adaptive collective representation**. One important direction is the development of **agent-aware shared memory**, where read and write behaviors are conditioned on agent roles, expertise, and trust, enabling more structured and reliable knowledge aggregation.

Another promising avenue lies in **learning-driven shared memory management**. Rather than relying on hand-designed policies for synchronization, summarization, or conflict resolution, future systems may train agents to decide when, what, and how to contribute to shared memory based on long-horizon team performance. Finally, as MAS increasingly operate in open-ended and multimodal environments, shared memory must support abstraction across heterogeneous signals while maintaining temporal and semantic coherence, for which we believe latent memory exhibits a promising path (Wu et al., 2025d). Advancing in these directions will be critical for scaling shared memory from a coordination aid into a foundation for robust collective intelligence.

## 7.6 Memory for World Model

### 7.6.1 Look-Back

The core objective of a World Model is to construct an internal environment capable of high-fidelity simulation of the physical world. These systems serve as the critical infrastructure for next-generation artificial intelligence. The core attribute of world model is to generate content that is both infinitely extensible and interactive in real time. Unlike traditional video generation that creates fixed-length clips, world models operate in an iterative manner by receiving actions at each step and predicting the next state to provide continuous feedback. In this iterative framework, the memory mechanism becomes the cornerstone of the system. Memory stores and maintains the spatial and semantic information or hidden states from the previous time step. It ensures that the generation of the next chunk maintains long-term consistency with the preceding context regarding scene layout, object attributes, and motion logic. Essentially, the memory mechanism enables world models to handle long-term temporal dependencies and realize trustworthy simulation interactions.

Previously, memory modeling relied on simplistic buffering approaches. Frame Sampling conditioned generation on a few historical frames (Bruce et al., 2024). While intuitive, this led to context fragmentation and perceptual drift as early details were lost. Sliding Window methods adapted LLM techniques like attention sinks and local KV caches (Liu et al., 2025e). Although this resolved computational bottlenecks, it restricted memory to a fixed window. Once an object left this view, the model effectively forgot it, preventing complex tasks like loop closure. By late 2025, the field shifted from finite context windows to structured state representations. Current architectures follow three main paths:

- State-Space Models (SSMs) architectures like Long-Context SSMs utilize Mamba-style backbones (Po et al., 2025; Yu et al., 2025f). These compress infinite history into a fixed-size recursive state, enabling theoretically infinite memory capacity with constant inference costs.

- Explicit Memory Banks. Unlike compressed states, these systems maintain an external storage of historical representations to support precise recall. Approaches differ in their structuring logic: UniWM employs a *hierarchical design*, separating short-term perception from long-term history via feature-based similarity gating (Dong et al., 2025b). Conversely, **retrieval-based approaches** like WorldMem and Context-as-Memory (CaM) maintain a flat bank of past contexts, utilizing *geometric retrieval* (e.g., FOV overlap) to dynamically select relevant frames for maintaining 3D scene consistency (Xiao et al., 2025c; Yu et al., 2025c).

- Sparse Memory and Retrieval To balance long-term adherence with efficiency, Genie Envisioner and Ctrl-World utilize sparse memory mechanisms (Liao et al., 2025b; Guo et al., 2025). These models augment current observations by injecting sparsely sampled historical frames or retrieving pose-conditioned context to anchor predictions and prevent drift during manipulation tasks.

### 7.6.2 Future Perspective

From an architectural perspective, the field is undergoing a fundamental transition from Data Caching which focuses on passive retention to State Simulation which focuses on active maintenance. This evolution is currently crystallizing into two distinct paradigms that aim to solve the conflict between real-time responsiveness and long-term logical consistency.

- The Dual-System Architecture. Inspired by cognitive science, world models could be bifurcated into fast and slow components. System 1 represents the fast and instinctive layer that handles immediate

physics and fluid interaction using efficient backbones like SSMs. System 2 represents the slow and deliberative layer that handles complex reasoning, planning, and world consistency using large-scale VLMs or explicit memory databases.

- **Active Memory Management.** Passive mechanisms are being superseded by Active Memory Policies. Instead of treating memory as a fixed buffer that blindly stores recent history, new models are designed as Cognitive Workspaces that actively curate, summarize, and discard information based on task relevance. Recent empirical studies demonstrate that such active memory management significantly outperforms static retrieval methods in handling functional infinite context. This shift marks the move from simply remembering the last N tokens to maintaining a coherent and queryable world state.

## 7.7 Trustworthy Memory

### 7.7.1 Look-Back: From Trustworthy RAG to Trustworthy Memory

As shown throughout this survey, memory plays a foundational role in enabling agentic behavior, which supports persistence, personalization, and continual learning. However, as memory systems become more deeply embedded into LLM-based agents, the question of *trustworthiness* has become paramount.

Earlier concerns around hallucination and factuality in retrieval-augmented generation (RAG) systems (Niu et al., 2024; Sun et al., 2025e; Lu et al., 2025c) have now evolved into a broader trust discourse for memory-augmented agents. Similar to RAG, one major motivation for using external or long-term memory is to reduce hallucinations by grounding model outputs in retrievable, factual content (Ru et al., 2024; Wang et al., 2025c). However, unlike RAG, agent memory often stores user-specific, persistent, and potentially sensitive content, ranging from factual knowledge to past interactions, preferences, or behavioral traces. This introduces additional challenges in privacy, interpretability, and safety.

Recent work by Wang et al. (2025b) demonstrates that memory modules can leak private data through indirect prompt-based attacks, highlighting the risk of memorization and over-retention. Concurrently, Wu et al. (2025g) argues that agent memory systems must support explicit mechanisms for *access control*, *verifiable forgetting*, and *auditable updates* to remain trustworthy. Notably, such threats are magnified in agent scenarios where memory persists across long time horizons.

Explainability also remains a critical bottleneck. While explicit memory, such as text logs or key-value stores, offers some transparency, users and developers still lack tools to trace which memory items were retrieved, how they influenced generation, or whether they were misused. In this regard, diagnostic tools like RAGChecker (Ru et al., 2024) and conflict-resolution frameworks such as RAMDocs with MADAM-RAG (Wang et al., 2025d) provide inspiration for tracing memory usage and reasoning under uncertainty.

Moreover, beyond individual memory, Shi et al. (2025d) and Rezazadeh et al. (2025a) highlight the emerging importance of *collective privacy* in shared or federated memory systems, which may operate across multi-agent deployments or organizations. All these developments collectively signal a need to elevate trust as a first-class principle in memory design.

### 7.7.2 Future Perspective

Looking ahead, we argue that **trustworthy memory** must be built around three interlinked pillars: *privacy preservation*, *explainability*, and *hallucination robustness*—each demanding architectural and algorithmic innovations.

For privacy, future systems should support granular permissioned memory, user-governed retention policies, encrypted or on-device storage, and federated access where needed (Wu et al., 2025g; Shi et al., 2025d; Rezazadeh et al., 2025a). Techniques like differential privacy, memory redaction, and adaptive forgetting (e.g., decay-based models or user-erasure interfaces) can serve as safeguards against both memorization and leakage (Chhikara et al., 2025).

Explainability requires moving beyond visible content to include *traceable access paths*, *self-rationalizing retrievals*, and possibly counterfactual reasoning (e.g., what would have changed without this memory?) (Ope-

nAI, 2024; Zhang et al., 2025u). Visualizations of memory attention, causal graphs of memory influence, and user-facing debugging tools may become standard components.

Hallucination mitigation will benefit from continued advances in conflict detection, multi-document reasoning, and uncertainty-aware generation. Strategies such as abstention under low-confidence retrieval, fallback to model priors (Wang et al., 2025c), or multi-agent cross-checking (Hu et al., 2024) are promising. Beyond behavioral safeguards, emerging *mechanistic interpretability* techniques offer a complementary direction by analyzing how internal representations and reasoning circuits contribute to hallucinated outputs. Methods such as representation-level probing and reasoning-path decomposition enable finer-grained diagnosis of where hallucinations originate, and provide principled tools for intervention and control (Sun et al., 2025e,c).

In the long term, we envision memory systems governed by OS-like abstractions: segmented, version-controlled, auditable, and jointly managed by agent and user (Packer et al., 2023b). Building such systems will require coordinated efforts across representation learning, system design, and policy control. As LLM agents begin to operate in persistent, open-ended environments, trustworthy memory will not just be a desirable feature—but a foundational requirement for real-world deployment.

## 7.8 Human–Cognitive Connections

### 7.8.1 Look Back

The architecture of contemporary agent memory systems has converged with foundational models of human cognition established over the last century. The prevailing design, which couples a capacity-limited context window with massive external vector databases, mirrors the Atkinson-Shiffrin multi-store model (Atkinson and Shiffrin, 1968), effectively instantiating an artificial counterpart to the distinction between working memory and long-term memory (Baddeley, 2012). Furthermore, the partitioning of agent memory into interaction logs, world knowledge, and code-based skills exhibits a striking structural alignment with *Tulving's* classification of *episodic*, *semantic*, and *procedural* memory (Tulving, 1972; Squire, 2004). Current frameworks (Zhong et al., 2024; Park et al., 2023; Gutierrez et al., 2024; Li et al., 2025k) operationalize these biological categories into engineering artifacts, where episodic memory provides autobiographical continuity and semantic memory offers generalized world knowledge.

Despite these structural parallels, a fundamental divergence remains in the *dynamics* of retrieval and maintenance. Human memory operates as a *constructive process*, where the brain actively reconstructs past events based on current cognitive states rather than replaying exact recordings (Schacter and Addis, 2007). In contrast, the majority of existing agent memory systems rely on verbatim retrieval mechanisms like RAG, treating memory as a repository of *immutable* tokens to be queried via semantic similarity (Packer et al., 2023b; Chhikara et al., 2025). Consequently, while agents possess a veridical record of the past, they lack the biological capacity for memory distortion, abstraction, and the dynamic remodeling of history that characterizes human intelligence.

### 7.8.2 Future Perspective

To bridge the gap between static storage and dynamic cognition, the next generation of agents must evolve beyond exclusive **online updating** by incorporating **offline consolidation** mechanisms analogous to biological sleep. Drawing from the Complementary Learning Systems (CLS) theory (Kumaran et al., 2016; McClelland et al., 1995), future architectures will likely introduce dedicated consolidation intervals where agents decouple from environmental interaction to engage in memory reorganization and generative replay (Mattar and Daw, 2018). During these offline periods, agents can autonomously distill generalizable schemas from raw episodic traces, perform **active forgetting** to prune redundant noise (Anderson and Hulbert, 2021), and optimize their internal indices without the latency constraints of real-time processing.

Ultimately, this evolution suggests a paradigm shift in memory forms and functions: moving from explicit text retrieval to **generative reconstruction**. Future systems may utilize generative memory (Zhang et al., 2025d) where the agent synthesizes latent memory tokens on demand, mirroring the brain's reconstructive nature. By integrating sleep-like consolidation cycles, agents will evolve from entities that merely archive data to those that internalize experience, resolving the stability-plasticity dilemma by periodically compacting vast episodic streams into efficient, parametric intuition.

# 8 Conclusion

This survey has examined agent memory as a foundational component of modern LLM-based agentic systems. By framing existing research through the unified lenses of *forms, functions, and dynamics*, we have clarified the conceptual landscape of agent memory and situated it within the broader evolution of agentic intelligence. On the level of **forms**, we identify three principal realizations: token-level, parametric, and latent memory, each of which has undergone distinct and rapid advances in recent years, reflecting fundamentally different trade-offs in representation, adaptability, and integration with agent policies. On the level of **functions**, we move beyond the coarse long-term versus short-term dichotomy prevalent in prior surveys, and instead propose a more fine-grained and encompassing taxonomy that distinguishes *factual, experiential, and working memory* according to their roles in knowledge retention, capability accumulation, and task-level reasoning. Together, these perspectives reveal that memory is not merely an auxiliary storage mechanism, but an essential substrate through which agents achieve temporal coherence, continual adaptation, and long-horizon competence.

Beyond organizing prior work, we have identified key challenges and emerging directions that point toward the next stage of agent memory research. In particular, the increasing integration of reinforcement learning, the rise of multimodal and multi-agent settings, and the shift from retrieval-centric to generative memory paradigms suggest a future in which memory systems become fully learnable, adaptive, and self-organizing. Such systems hold the potential to transform large language models from powerful but static generators into agents capable of sustained interaction, self-improvement, and principled reasoning over time.

We hope this survey provides a coherent foundation for future research and serves as a reference for both researchers and practitioners. As agentic systems continue to mature, the design of memory will remain a central and open problem, one that is likely to play a decisive role in the development of robust, general, and enduring artificial intelligence.

# References

Aadharsh Aadhithya A, Sachin Kumar S, and Soman K. P. Enhancing long-term memory using hierarchical aggregate tree for retrieval augmented generation, 2024. https://arxiv.org/abs/2406.06124.

Saaket Agashe, Jiuzhou Han, Shuyu Gan, Jiachen Yang, Ang Li, and Xin Eric Wang. Agent S: An Open Agentic Framework that Uses Computers Like a Human. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net, 2025.

AgentScope. GitHub - agentscope-ai/ReMe: ReMe: Memory Management Kit for Agents - Remember Me, Refine Me. — github.com. https://github.com/agentscope-ai/ReMe, 2025. [Accessed 14-12-2025].

Qingyao Ai, Yichen Tang, Changyue Wang, Jianming Long, Weihang Su, and Yiqun Liu. Memorybench: A benchmark for memory and continual learning in llm systems. *arXiv preprint arXiv:2510.17281*, 2025.

Reem Aleithan, Haoran Xue, Mohammad Mahdi Mohajer, Elijah Nnorom, Gias Uddin, and Song Wang. Swe-bench+: Enhanced coding benchmark for llms, 2024. https://arxiv.org/abs/2410.06992.

Nick Alonso, Tomas Figliolia, Anthony Ndirango, and Beren Millidge. Toward conversational agents with context and time sensitive long-term memory. *CoRR*, abs/2406.00057, 2024. doi: 10.48550/ARXIV.2406.00057. https://doi.org/10.48550/arXiv.2406.00057.

Michael C. Anderson and Justin C. Hulbert. Active forgetting: Adaptation of memory by prefrontal control. *Annual Review of Psychology*, 72:1–36, January 2021. ISSN 1545-2085. doi: 10.1146/annurev-psych-072720-094140.

Petr Anokhin, Nikita Semenov, Artyom Sorokin, Dmitry Evseev, Andrey Kravchenko, Mikhail Burtsev, and Evgeny Burnaev. Arigraph: Learning knowledge graph world models with episodic memory for llm agents. *arXiv preprint arXiv:2407.04363*, 2024.

Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to retrieve, generate, and critique through self-reflection, 2023. https://arxiv.org/abs/2310.11511.

R. C. Atkinson and R. M. Shiffrin. Human memory: A proposed system and its control processes. In *The Psychology of Learning and Motivation: II*, pages xi, 249–xi, 249. Academic Press, Oxford, England, 1968. doi: 10.1016/S0079-7421(08)60422-3.

Alan Baddeley. Working Memory: Theories, Models, and Controversies. *Annual Review of Psychology*, 63(Volume 63, 2012):1–29, January 2012. ISSN 0066-4308, 1545-2085. doi: 10.1146/annurev-psych-120710-100422.

Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. LongBench: A bilingual, multitask benchmark for long context understanding. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3119–3137, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.172. https://aclanthology.org/2024.acl-long.172/.

Yushi Bai, Shangqing Tu, Jiajie Zhang, Hao Peng, Xiaozhi Wang, Xin Lv, Shulin Cao, Jiazheng Xu, Lei Hou, Yuxiao Dong, et al. Longbench v2: Towards deeper understanding and reasoning on realistic long-context multitasks. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3639–3664, 2025.

Alexandre Bailly, Antoine Saubin, Gabriel Kocevar, and Jonathan Bodin. Divide and summarize: improve slm text summarization. *Frontiers in Artificial Intelligence*, 8:1604034, 2025.

Ali Behrouz, Meisam Razaviyayn, Peilin Zhong, and Vahab Mirrokni. Nested learning: The illusion of deep learning architectures. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025a. https://openreview.net/forum?id=nbMeRvNb7A.

Ali Behrouz, Peilin Zhong, and Vahab Mirrokni. Titans: Learning to memorize at test time. *CoRR*, abs/2501.00663, 2025b. doi: 10.48550/ARXIV.2501.00663. https://doi.org/10.48550/arXiv.2501.00663.

Weihao Bo, Shan Zhang, Yanpeng Sun, Jingjing Wu, Qunyi Xie, Xiao Tan, Kunbin Chen, Wei He, Xiaofan Li, Na Zhao, Jingdong Wang, and Zechao Li. Agentic learner with grow-and-refine multimodal semantic memory, 2025. https://arxiv.org/abs/2511.21678.

Islem Bouzenia, Premkumar Devanbu, and Michael Pradel. RepairAgent: An Autonomous, LLM-Based Agent for Program Repair, October 2024. http://arxiv.org/abs/2403.17134. arXiv:2403.17134 [cs].

Jake Bruce, Michael D Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. Genie: Generative interactive environments. In *Forty-first International Conference on Machine Learning*, 2024.

Yuzheng Cai, Siqi Cai, Yuchen Shi, Zihan Xu, Lichao Chen, Yulei Qin, Xiaoyu Tan, Gang Li, Zongyi Li, Haojia Lin, Yong Mao, Ke Li, and Xing Sun. Training-free group relative policy optimization, 2025a. https://arxiv.org/abs/2510.08191.

Zefan Cai, Yichi Zhang, Bofei Gao, Yuliang Liu, Tianyu Liu, Keming Lu, Wayne Xiong, Yue Dong, Baobao Chang, Junjie Hu, and Wen Xiao. Pyramidkv: Dynamic KV cache compression based on pyramidal information funneling. *CoRR*, abs/2406.02069, 2024. doi: 10.48550/ARXIV.2406.02069. https://doi.org/10.48550/arXiv.2406.02069.

Zhicheng Cai, Xinyuan Guo, Yu Pei, Jiangtao Feng, Jinsong Su, Jiangjie Chen, Ya-Qin Zhang, Wei-Ying Ma, Mingxuan Wang, and Hao Zhou. Flex: Continuous agent evolution via forward learning from experience, 2025b. https://arxiv.org/abs/2511.06449.

CAMEL-AI. Workforce — camel-ai documentation. https://docs.camel-ai.org/key_modules/workforce, 2025. Accessed: 2025-08-09.

Nicola De Cao, Wilker Aziz, and Ivan Titov. Editing factual knowledge in language models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6491–6506, 2021.

Maitreyi Chatterjee and Devansh Agarwal. Semantic anchoring in agentic memory: Leveraging linguistic structures for persistent conversational context. *CoRR*, abs/2508.12630, 2025. doi: 10.48550/ARXIV.2508.12630. https://doi.org/10.48550/arXiv.2508.12630.

Baian Chen, Chang Shu, Ehsan Shareghi, Nigel Collier, Karthik Narasimhan, and Shunyu Yao. Fireact: Toward language agent fine-tuning, 2023a. https://arxiv.org/abs/2310.05915.

Dake Chen, Hanbin Wang, Yunhao Huo, Yuzhao Li, and Haoyang Zhang. GameGPT: Multi-agent Collaborative Framework for Game Development. *CoRR*, abs/2310.08067, 2023b. doi: 10.48550/ARXIV.2310.08067.

Guoxin Chen, Zile Qiao, Xuanzhong Chen, Donglei Yu, Haotian Xu, Wayne Xin Zhao, Ruihua Song, Wenbiao Yin, Huifeng Yin, Liwen Zhang, Kuan Li, Minpeng Liao, Yong Jiang, Pengjun Xie, Fei Huang, and Jingren Zhou. Iterresearch: Rethinking long-horizon agents via markovian state reconstruction, 2025a. https://arxiv.org/abs/2511.07327.

Kaiyuan Chen, Yixin Ren, Yang Liu, Xiaobo Hu, Haotong Tian, Tianbao Xie, Fangfu Liu, Haoye Zhang, Hongzhang Liu, Yuan Gong, Chen Sun, Han Hou, Hui Yang, James Pan, Jianan Lou, Jiayi Mao, Jizheng Liu, Jinpeng Li, Kangyi Liu, Kenkun Liu, Rui Wang, Run Li, Tong Niu, Wenlong Zhang, Wenqi Yan, Xuanzheng Wang, Yuchen Zhang, Yi-Hsin Hung, Yuan Jiang, Zexuan Liu, Zihan Yin, Zijian Ma, and Zhiwen Mo. xbench: Tracking agents productivity scaling with profession-aligned real-world evaluations, 2025b. https://arxiv.org/abs/2506.13651.

Nuo Chen, Hongguang Li, Jianhui Chang, Juhua Huang, Baoyuan Wang, and Jia Li. Compress to impress: Unleashing the potential of compressive memory in real-world long-term conversations. In Owen Rambow, Leo Wanner, Marianna Apidianaki, Hend Al-Khalifa, Barbara Di Eugenio, and Steven Schockaert, editors, *Proceedings of the 31st International Conference on Computational Linguistics, COLING 2025, Abu Dhabi, UAE, January 19-24, 2025*, pages 755–773. Association for Computational Linguistics, 2025c. https://aclanthology.org/2025.coling-main.51/.

Qiuhui Chen, Qiang Fu, Hao Bai, and Yi Hong. Longformer: Longitudinal transformer for alzheimer's disease classification with structural mris. In *IEEE/CVF Winter Conference on Applications of Computer Vision, WACV 2024, Waikoloa, HI, USA, January 3-8, 2024*, pages 3563–3572. IEEE, 2024a. doi: 10.1109/WACV57701.2024.00354. https://doi.org/10.1109/WACV57701.2024.00354.

Weishu Chen, Jinyi Tang, Zhouhui Hou, Shihao Han, Mingjie Zhan, Zhiyuan Huang, Delong Liu, Jiawei Guo, Zhicheng Zhao, and Fei Su. Moom: Maintenance, organization and optimization of memory in ultra-long role-playing dialogues, 2025d. https://arxiv.org/abs/2509.11860.

Xiuying Chen, Shen Gao, Mingzhe Li, Qingqing Zhu, Xin Gao, and Xiangliang Zhang. Write summary step-by-step: A pilot study of stepwise summarization. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 32: 1406–1415, 2024b.

Yinpeng Chen, DeLesley Hutchins, Aren Jansen, Andrey Zhmoginov, David Racz, and Jesper Andersen. Melodi: Exploring memory compression for long contexts, 2024c. https://arxiv.org/abs/2410.03156.

Zhaorun Chen, Zhuokai Zhao, Kai Zhang, Bo Liu, Qi Qi, Yifan Wu, Tarun Kalluri, Sara Cao, Yuanhao Xiong, Haibo Tong, Huaxiu Yao, Hengduo Li, Jiacheng Zhu, Xian Li, Dawn Song, Bo Li, Jason Weston, and Dat Huynh. Scaling agent learning via experience synthesis, 2025e. https://arxiv.org/abs/2511.03773.

Ho Kei Cheng and Alexander G. Schwing. Xmem: Long-term video object segmentation with an atkinson-shiffrin memory model, 2022. https://arxiv.org/abs/2207.07115.

Alexis Chevalier, Alexander Wettig, Anirudh Ajith, and Danqi Chen. Adapting language models to compress contexts. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 3829–3846. Association for Computational Linguistics, 2023. doi: 10.18653/V1/2023.EMNLP-MAIN.232. https://doi.org/10.18653/v1/2023.emnlp-main.232.

Maxime Chevalier-Boisvert, Dzmitry Bahdanau, Salem Lahlou, Lucas Willems, Chitwan Saharia, Thien Huu Nguyen, and Yoshua Bengio. Babyai: A platform to study the sample efficiency of grounded language learning, 2019. https://arxiv.org/abs/1810.08272.

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

Eunseong Choi, June Park, Hyeri Lee, and Jongwuk Lee. Conflict-aware soft prompting for retrieval-augmented generation. *CoRR*, abs/2508.15253, 2025. doi: 10.48550/ARXIV.2508.15253. https://doi.org/10.48550/arXiv.2508.15253.

Cognee. GitHub - topoteretes/cognee: Memory for AI Agents in 6 lines of code. https://github.com/topoteretes/cognee, 2025. [Accessed 14-12-2025].

Nelson Cowan. Working Memory Underpins Cognitive Development, Learning, and Education. *Educational psychology review*, 26(2):197–223, June 2014. ISSN 1040-726X. doi: 10.1007/s10648-013-9246-y.

Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. https://openreview.net/forum?id=mZn2Xyh9Ec.

Nicola De Cao, Wilker Aziz, and Ivan Titov. Editing factual knowledge in language models. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6491–6506, Online and Punta Cana, Dominican

Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.522. https://aclanthology.org/2021.emnlp-main.522/.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. https://arxiv.org/abs/2501.12948.

Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samual Stevens, Boshi Wang, Huan Sun, and Yu Su. Mind2Web: Towards a Generalist Agent for the Web. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. http://papers.nips.cc/paper_files/paper/2023/hash/5950bf290a1570ea401bf98882128160-Abstract-Datasets_and_Benchmarks.html.

Yang Deng, Xuan Zhang, Wenxuan Zhang, Yifei Yuan, See-Kiong Ng, and Tat-Seng Chua. On the multi-turn instruction following for conversational web agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8795–8812, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.477. https://aclanthology.org/2024.acl-long.477/.

Junnan Dong, Siyu An, Yifei Yu, Qian-Wen Zhang, Linhao Luo, Xiao Huang, Yunsheng Wu, Di Yin, and Xing Sun. Youtu-graphrag: Vertically unified agents for graph retrieval-augmented complex reasoning, 2025a. https://arxiv.org/abs/2508.19855.

Yifei Dong, Fengyi Wu, Guangyu Chen, Zhi-Qi Cheng, Qiyu Hu, Yuxuan Zhou, Jingdong Sun, Jun-Yan He, Qi Dai, and Alexander G Hauptmann. Unified world models: Memory-augmented planning and foresight for visual navigation. *arXiv preprint arXiv:2510.08713*, 2025b.

Yiming Du, Hongru Wang, Zhengyi Zhao, Bin Liang, Baojun Wang, Wanjun Zhong, Zezhong Wang, and Kam-Fai Wong. PerLTQA: A personal long-term memory dataset for memory classification, retrieval, and fusion in question answering. In *Proceedings of the 10th SIGHAN Workshop on Chinese Language Processing (SIGHAN-10)*, pages 152–164, Bangkok, Thailand, August 2024. Association for Computational Linguistics. https://aclanthology.org/2024.sighan-1.18/.

Yiming Du, Wenyu Huang, Danna Zheng, Zhaowei Wang, Sebastien Montella, Mirella Lapata, Kam-Fai Wong, and Jeff Z Pan. Rethinking memory in ai: Taxonomy, operations, topics, and future directions. *arXiv preprint arXiv:2505.00675*, 2025a.

Yiming Du, Bingbing Wang, Yang He, Bin Liang, Baojun Wang, Zhongyang Li, Lin Gui, Jeff Z. Pan, Ruifeng Xu, and Kam-Fai Wong. Memguide: Intent-driven memory selection for goal-oriented multi-session llm agents, 2025b. https://arxiv.org/abs/2505.20231.

Zane Durante, Qiuyuan Huang, Naoki Wake, Ran Gong, Jae Sung Park, Bidipta Sarkar, Rohan Taori, Yusuke Noda,

Demetri Terzopoulos, Yejin Choi, Katsushi Ikeuchi, Hoi Vo, Li Fei-Fei, and Jianfeng Gao. Agent ai: Surveying the horizons of multimodal interaction, 2024. https://arxiv.org/abs/2401.03568.

Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, Dasha Metropolitansky, Robert Osazuwa Ness, and Jonathan Larson. From local to global: A graph rag approach to query-focused summarization, 2025. https://arxiv.org/abs/2404.16130.

Yue Fan, Xiaojian Ma, Rongpeng Su, Jun Guo, Rujie Wu, Xi Chen, and Qing Li. Embodied videoagent: Persistent memory from egocentric videos and embodied sensors enables dynamic scene understanding. *CoRR*, abs/2501.00358, 2025. doi: 10.48550/ARXIV.2501.00358. https://doi.org/10.48550/arXiv.2501.00358.

Jinyuan Fang, Yanwen Peng, Xi Zhang, Yingxu Wang, Xinhao Yi, Guibin Zhang, Yi Xu, Bin Wu, Siwei Liu, Zihao Li, Zhaochun Ren, Nikos Aletras, Xi Wang, Han Zhou, and Zaiqiao Meng. A comprehensive survey of self-evolving ai agents: A new paradigm bridging foundation models and lifelong agentic systems, 2025a. https://arxiv.org/abs/2508.07407.

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*, 2025b.

Junfeng Fang, Houcheng Jiang, Kun Wang, Yunshan Ma, Shi Jie, Xiang Wang, Xiangnan He, and Tat seng Chua. Alphaedit: Null-space constrained knowledge editing for language models, 2025c. https://arxiv.org/abs/2410.02355.

Runnan Fang, Yuan Liang, Xiaobin Wang, Jialong Wu, Shuofei Qiao, Pengjun Xie, Fei Huang, Huajun Chen, and Ningyu Zhang. Memp: Exploring agent procedural memory, 2025d. https://arxiv.org/abs/2508.06433.

Erhu Feng, Wenbo Zhou, Zibin Liu, Le Chen, Yunpeng Dong, Cheng Zhang, Yisheng Zhao, Dong Du, Zhi-Hua Zhou, Yubin Xia, and Haibo Chen. Get experience from practice: LLM agents with record & replay. *CoRR*, abs/2505.17716, 2025. doi: 10.48550/ARXIV.2505.17716. https://doi.org/10.48550/arXiv.2505.17716.

Zafeirios Fountas, Martin Benfeghoul, Adnan Oomerjee, Fenia Christopoulou, Gerasimos Lampouras, Haitham Bou-Ammar, and Jun Wang. Human-inspired episodic memory for infinite context llms. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net, 2025. https://openreview.net/forum?id=BI2int5SAC.

Chen Gao, Xiaochong Lan, Zhihong Lu, Jinzhu Mao, Jinghua Piao, Huandong Wang, Depeng Jin, and Yong Li. S\(^\mbox3\): Social-network Simulation System with Large Language Model-Empowered Agents. *CoRR*, abs/2307.14984, 2023a. doi: 10.48550/ARXIV.2307.14984.

Hang Gao and Yongfeng Zhang. Memory sharing for large language model based agents. *CoRR*, abs/2404.09982, 2024a. doi: 10.48550/ARXIV.2404.09982. https://doi.org/10.48550/arXiv.2404.09982.

Hang Gao and Yongfeng Zhang. Memory Sharing for Large Language Model based Agents. *CoRR*, abs/2404.09982, 2024b. doi: 10.48550/ARXIV.2404.09982.

Hang Gao and Yongfeng Zhang. PTR: Precision-Driven Tool Recommendation for Large Language Models. *CoRR*, abs/2411.09613, 2024c. doi: 10.48550/ARXIV.2411.09613. https://doi.org/10.48550/arXiv.2411.09613. arXiv: 2411.09613.

Huan-ang Gao, Jiayi Geng, Wenyue Hua, Mengkang Hu, Xinzhe Juan, Hongzhang Liu, Shilong Liu, Jiahao Qiu, Xuan Qi, Yiran Wu, Hongru Wang, Han Xiao, Yuhang Zhou, Shaokun Zhang, Jiayi Zhang, Jinyu Xiang, Yixiong Fang, Qiwen Zhao, Dongrui Liu, Qihan Ren, Cheng Qian, Zhenhailong Wang, Minda Hu, Huazheng Wang, Qingyun Wu, Heng Ji, and Mengdi Wang. A survey of self-evolving agents: On path to artificial super intelligence, August 2025.

Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. Precise zero-shot dense retrieval without relevance labels. In Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 1762–1777. Association for Computational Linguistics, 2023b. doi: 10.18653/V1/2023.ACL-LONG.99. https://doi.org/10.18653/v1/2023.acl-long.99.

Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. Retrieval-augmented generation for large language models: A survey, 2024. https://arxiv.org/abs/2312.10997.

Tao Ge, Jing Hu, Lei Wang, Xun Wang, Si-Qing Chen, and Furu Wei. In-context Autoencoder for Context Compression in a Large Language Model. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.

Samuel J. Gershman, Ila Fiete, and Kazuki Irie. Key-value memory in the brain. *CoRR*, abs/2501.02950, 2025. doi: 10.48550/ARXIV.2501.02950. https://doi.org/10.48550/arXiv.2501.02950.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces, 2024. https://arxiv.org/abs/2312.00752.

Yanjiang Guo, Lucy Xiaoyang Shi, Jianyu Chen, and Chelsea Finn. Ctrl-world: A controllable generative world model for robot manipulation. *arXiv preprint arXiv:2510.10125*, 2025.

Saket Gurukar and Asim Kadav. Long-vmnet: Accelerating long-form video understanding via fixed memory, 2025. https://arxiv.org/abs/2503.13707.

Bernal Jimenez Gutierrez, Yiheng Shu, Yu Gu, Michihiro Yasunaga, and Yu Su. HippoRAG: Neurobiologically inspired long-term memory for large language models. In *Advances in Neural Information Processing Systems*, 2024.

Bernal Jiménez Gutiérrez, Yiheng Shu, Weijian Qi, Sizhe Zhou, and Yu Su. From rag to memory: Non-parametric continual learning for large language models, 2025. https://arxiv.org/abs/2502.14802.

Dongge Han, Camille Couturier, Daniel Madrigal Diaz, Xuchao Zhang, Victor Rühle, and Saravan Rajmohan. LEGOMem: Modular Procedural Memory for Multi-agent LLM Systems for Workflow Automation, October 2025a. http://arxiv.org/abs/2510.04851. arXiv:2510.04851 [cs].

Haoyu Han, Yu Wang, Harry Shomer, Kai Guo, Jiayuan Ding, Yongjia Lei, Mahantesh Halappanavar, Ryan A. Rossi, Subhabrata Mukherjee, Xianfeng Tang, Qi He, Zhigang Hua, Bo Long, Tong Zhao, Neil Shah, Amin Javari, Yinglong Xia, and Jiliang Tang. Retrieval-augmented generation with graphs (graphrag), 2025b. https://arxiv.org/abs/2501.00309.

Jinyi Han, Xinyi Wang, Haiquan Zhao, Tingyun li, Zishang Jiang, Sihang Jiang, Jiaqing Liang, Xin Lin, Weikang Zhou, Zeye Sun, Fei Yu, and Yanghua Xiao. A stitch in time saves nine: Proactive self-refinement for language models, 2025c. https://arxiv.org/abs/2508.12903.

Jackson Hassell, Dan Zhang, Hannah Kim, Tom Mitchell, and Estevam Hruschka. Learning from supervision with semantic and episodic memory: A reflective approach to agent adaptation. *arXiv preprint arXiv:2510.19897*, 2025.

Bo He, Hengduo Li, Young Kyun Jang, Menglin Jia, Xuefei Cao, Ashish Shah, Abhinav Shrivastava, and Ser-Nam Lim. MA-LMM: Memory-Augmented Large Multimodal Model for Long-Term Video Understanding. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*, pages 13504–13514. IEEE, 2024. doi: 10.1109/CVPR52733.2024.01282.

Junqing He, Liang Zhu, Rui Wang, Xi Wang, Gholamreza Haffari, and Jiaxing Zhang. MADial-bench: Towards real-world evaluation of memory-augmented dialogue generation. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 9902–9921, Albuquerque, New Mexico, April 2025. Association for Computational Linguistics. ISBN 979-8-89176-189-6. doi: 10.18653/v1/2025.naacl-long.499. https://aclanthology.org/2025.naacl-long.499/.

Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: Decoding-enhanced bert with disentangled attention. *arXiv preprint arXiv:2006.03654*, 2020.

Dan Hendrycks, Dawn Song, Christian Szegedy, Honglak Lee, Yarin Gal, Erik Brynjolfsson, Sharon Li, Andy Zou, Lionel Levine, Bo Han, et al. A definition of agi. *arXiv preprint arXiv:2510.18212*, 2025.

Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning steps, 2020. https://arxiv.org/abs/2011.01060.

Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. MetaGPT: Meta programming for A multi-agent collaborative framework. In *The Twelfth International Conference on Learning Representations*, 2024.

Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekesh, Fei Jia, Yang Zhang, and Boris Ginsburg. Ruler: What's the real context size of your long-context language models? *arXiv preprint arXiv:2404.06654*, 2024.

Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *The Tenth International Conference on Learning Representations*, 2022.

Mengkang Hu, Tianxing Chen, Qiguang Chen, Yao Mu, Wenqi Shao, and Ping Luo. HiAgent: Hierarchical Working Memory Management for Solving Long-Horizon Agent Tasks with Large Language Model. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar, editors, *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pages 32779–32798. Association for Computational Linguistics, 2025a.

Shengran Hu, Cong Lu, and Jeff Clune. Automated design of agentic systems. In *The Thirteenth International Conference on Learning Representations*, 2025b.

Xiangkun Hu, Dongyu Ru, Lin Qiu, Qipeng Guo, Tianhang Zhang, Yang Xu, Yun Luo, Pengfei Liu, Yue Zhang, and Zheng Zhang. Refchecker: Reference-based fine-grained hallucination checker and benchmark for large language models. *arXiv preprint arXiv:2405.14486*, 2024.

Yuanzhe Hu, Yu Wang, and Julian McAuley. Evaluating memory in llm agents via incremental multi-turn interactions. *arXiv preprint arXiv:2507.05257*, 2025c.

Jen-tse Huang, Kaiser Sun, Wenxuan Wang, and Mark Dredze. Language Models Do Not Have Human-Like Working Memory, September 2025a.

Jiani Huang, Xingchen Zou, Lianghao Xia, and Qing Li. Mr.rec: Synergizing memory and reasoning for personalized recommendation assistant with llms. *CoRR*, abs/2510.14629, 2025b. doi: 10.48550/ARXIV.2510.14629. https://doi.org/10.48550/arXiv.2510.14629.

Tenghao Huang, Kinjal Basu, Ibrahim Abdelaziz, Pavan Kapanipathi, Jonathan May, and Muhao Chen. R2D2: Remembering, Replaying and Dynamic Decision Making with a Reflective Agentic Memory. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar, editors, *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pages 30318–30330. Association for Computational Linguistics, 2025c. https://aclanthology.org/2025.acl-long.1464/.

Xu Huang, Jianxun Lian, Yuxuan Lei, Jing Yao, Defu Lian, and Xing Xie. Recommender AI agent: Integrating large language models for interactive recommendations. *ACM Trans. Inf. Syst.*, 43(4):96:1–96:33, 2025d. doi: 10.1145/3731446. https://doi.org/10.1145/3731446.

Zeyu Huang, Yikang Shen, Xiaofeng Zhang, Jie Zhou, Wenge Rong, and Zhang Xiong. Transformer-patcher: One mistake worth one neuron, 2023. https://arxiv.org/abs/2301.09785.

Md Ashraful Islam, Mohammed Eunus Ali, and Md Rizwan Parvez. MapCoder: Multi-Agent Code Generation for Competitive Problem Solving. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 4912–4944. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.ACL-LONG.269. https://doi.org/10.18653/v1/2024.acl-long.269.

Kai Tzu iunn Ong, Namyoung Kim, Minju Gwak, Hyungjoo Chae, Taeyoon Kwon, Yohan Jo, Seung won Hwang, Dongha Lee, and Jinyoung Yeo. Towards lifelong dialogue agents via timeline-based memory management, 2025. https://arxiv.org/abs/2406.10996.

Jingyi Jia and Qinbin Li. Autotool: Efficient tool selection for large language model agents, November 2025.

Zixi Jia, Qinghua Liu, Hexiao Li, Yuyan Chen, and Jiqiang Liu. Evaluating the long-term memory of large language models. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 19759–19777, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-256-5. doi: 10.18653/v1/2025.findings-acl.1014. https://aclanthology.org/2025.findings-acl.1014/.

Bowen Jiang, Zhuoqun Hao, Young-Min Cho, Bryan Li, Yuan Yuan, Sihao Chen, Lyle Ungar, Camillo J Taylor, and Dan Roth. Know me, respond to me: Benchmarking llms for dynamic user profiling and personalized responses at scale. *arXiv preprint arXiv:2504.14225*, 2025a.

Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. LLMLingua: Compressing Prompts for Accelerated Inference of Large Language Models. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore,*

*December 6-10, 2023*, pages 13358–13376. Association for Computational Linguistics, 2023. doi: 10.18653/V1/2023. EMNLP-MAIN.825.

Huiqiang Jiang, Qianhui Wu, Xufang Luo, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. LongLLMLingua: Accelerating and Enhancing LLMs in Long Context Scenarios via Prompt Compression. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 1658–1677. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.ACL-LONG.91.

Jiantong Jiang, Peiyu Yang, Rui Zhang, and Feng Liu. Towards efficient large language model serving: A survey on system-aware KV cache optimization. *TechRxiv*, 2025b. doi: 10.36227/techrxiv.176046306.66521015/v2. http://dx.doi.org/10.36227/techrxiv.176046306.66521015/v2.

Tao Jiang, Zichuan Lin, Lihe Li, Yi-Chen Li, Cong Guan, Lei Yuan, Zongzhang Zhang, Yang Yu, and Deheng Ye. Multi-agent in-context coordination via decentralized memory retrieval. *arXiv preprint arXiv:2511.10030*, 2025c.

Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R. Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*, 2024.

Bowen Jin, Hansi Zeng, Zhenrui Yue, Dong Wang, Hamed Zamani, and Jiawei Han. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *CoRR*, abs/2503.09516, 2025. doi: 10.48550/ARXIV.2503. 09516. https://doi.org/10.48550/arXiv.2503.09516.

Kumara Kahatapitiya, Kanchana Ranasinghe, Jongwoo Park, and Michael S Ryoo. Language repository for long video understanding. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 5627–5646, 2025.

Zhao Kaiya, Michelangelo Naim, Jovana Kondic, Manuel Cortes, Jiaxin Ge, Shuying Luo, Guangyu Robert Yang, and Andrew Ahn. Lyfe agents: Generative agents for low-cost real-time social interactions, 2023. https://arxiv.org/abs/2310.02172.

Jiazheng Kang, Mingming Ji, Zhe Zhao, and Ting Bai. Memory os of ai agent, 2025a. https://arxiv.org/abs/2506.06326.

Jikun Kang, Wenqi Wu, Filippos Christianos, Alex J. Chan, Fraser Greenlee, George Thomas, Marvin Purtorab, and Andy Toulis. LM2: large memory models. *CoRR*, abs/2502.06049, 2025b. doi: 10.48550/ARXIV.2502.06049. https://doi.org/10.48550/arXiv.2502.06049.

Minki Kang, Wei-Ning Chen, Dongge Han, Huseyin A. Inan, Lukas Wutschitz, Yanzhi Chen, Robert Sim, and Saravan Rajmohan. Acon: Optimizing context compression for long-horizon llm agents, October 2025c.

Dongkyu Kim, Byoungwook Kim, Donggeon Han, and Matous Eibich. Autorag: Automated framework for optimization of retrieval augmented generation pipeline. *CoRR*, abs/2410.20878, 2024a. doi: 10.48550/ARXIV.2410.20878. https://doi.org/10.48550/arXiv.2410.20878.

Gangwoo Kim, Sungdong Kim, Byeongguk Jeon, Joonsuk Park, and Jaewoo Kang. Tree of clarifications: Answering ambiguous questions with retrieval-augmented large language models, 2023a. https://arxiv.org/abs/2310.14696.

Hana Kim, Kai Tzu iunn Ong, Seoyeon Kim, Dongha Lee, and Jinyoung Yeo. Commonsense-augmented memory construction and management in long-term conversations via context-aware persona refinement, 2024b. https://arxiv.org/abs/2401.14215.

Hyuntak Kim and Byung-Hak Kim. Nexussum: Hierarchical llm agents for long-form narrative summarization. *arXiv preprint arXiv:2505.24575*, 2025.

Namyoung Kim, Kai Tzu-iunn Ong, Yeonjun Hwang, Minseok Kang, Iiseo Jihn, Gayoung Kim, Minju Kim, and Jinyoung Yeo. PRINCIPLES: synthetic strategy memory for proactive dialogue agents. *CoRR*, abs/2509.17459, 2025a. doi: 10.48550/ARXIV.2509.17459. https://doi.org/10.48550/arXiv.2509.17459.

Sangyeop Kim, Yohan Lee, Sanghwa Kim, Hyunjong Kim, and Sungzoon Cho. Pre-storage reasoning for episodic memory: Shifting inference burden to memory for personalized dialogue. *CoRR*, abs/2509.10852, 2025b. doi: 10.48550/ARXIV.2509.10852. https://doi.org/10.48550/arXiv.2509.10852.

Taewoon Kim, Michael Cochez, Vincent François-Lavet, Mark A. Neerincx, and Piek Vossen. A machine with short-term, episodic, and semantic memory systems. In Brian Williams, Yiling Chen, and Jennifer Neville, editors, *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial*

*Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023*, pages 48–56. AAAI Press, 2023b. doi: 10.1609/AAAI.V37I1.25075. https://doi.org/10.1609/aaai.v37i1.25075.

Dharshan Kumaran, Demis Hassabis, and James L. McClelland. What learning systems do intelligent agents need? complementary learning systems theory updated. *Trends in Cognitive Sciences*, 20(7):512–534, July 2016. ISSN 1879-307X. doi: 10.1016/j.tics.2016.05.004.

Yuri Kuratov, Aydar Bulatov, Petr Anokhin, Ivan Rodkin, Dmitry Sorokin, Artyom Y. Sorokin, and Mikhail Burtsev. Babilong: Testing the limits of llms with long context reasoning-in-a-haystack. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang, editors, *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024. http://papers.nips.cc/paper_files/paper/2024/hash/c0d62e70dbc659cc9bd44cbcf1cb652f-Abstract-Datasets_and_Benchmarks_Track.html.

Taeyoon Kwon, Dongwook Choi, Sunghwan Kim, Hyojun Kim, Seungjun Moon, Beong-woo Kwak, Kuan-Hao Huang, and Jinyoung Yeo. Embodied agents meet personalization: Exploring memory utilization for personalized assistance. *CoRR*, abs/2505.16348, 2025. doi: 10.48550/ARXIV.2505.16348. https://doi.org/10.48550/arXiv.2505.16348.

LangChain. GitHub - langchain-ai/langmem — github.com. https://github.com/langchain-ai/langmem, 2025. [Accessed 14-12-2025].

Gibbeum Lee, Volker Hartmann, Jongho Park, Dimitris Papailiopoulos, and Kangwook Lee. Prompted llms as chatbot modules for long open-domain conversation. In Anna Rogers, Jordan L. Boyd-Graber, and Naoaki Okazaki, editors, *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 4536–4554. Association for Computational Linguistics, 2023. doi: 10.18653/V1/2023.FINDINGS-ACL.277. https://doi.org/10.18653/v1/2023.findings-acl.277.

Kuang-Huei Lee, Xinyun Chen, Hiroki Furuta, John F. Canny, and Ian Fischer. A human-inspired reading agent with gist memory of very long contexts. In *Forty-first International Conference on Machine Learning*, 2024a.

Myeonghwa Lee, Seonho An, and Min-Soo Kim. Planrag: A plan-then-retrieval augmented generation for generative large language models as decision makers, 2024b. https://arxiv.org/abs/2406.12430.

Xiang Lei, Qin Li, and Min Zhang. D-smart: Enhancing llm dialogue consistency via dynamic structured memory and reasoning tree. *arXiv preprint arXiv:2510.13363*, 2025.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive NLP tasks. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html.

Baiqi Li, Zhiqiu Lin, Deepak Pathak, Jiayao Li, Yixin Fei, Kewen Wu, Tiffany Ling, Xide Xia, Pengchuan Zhang, Graham Neubig, and Deva Ramanan. Genai-bench: Evaluating and improving compositional text-to-visual generation, 2024a. https://arxiv.org/abs/2406.13743.

Caorui Li, Yu Chen, Yiyan Ji, Jin Xu, Zhenyu Cui, Shihao Li, Yuanxing Zhang, Jiafu Tang, Zhenghao Song, Dingling Zhang, Ying He, Haoxiang Liu, Yuxuan Wang, Qiufeng Wang, Zhenhe Wu, Jiehui Luo, Zhiyu Pan, Weihao Xie, Chenchen Zhang, Zhaohui Wang, Jiayi Tian, Yanghai Wang, Zhe Cao, Minxin Dai, Ke Wang, Runzhe Wen, Yinghao Ma, Yaning Pan, Sungkyun Chang, Termeh Taheri, Haiwen Xia, Christos Plachouras, Emmanouil Benetos, Yizhi Li, Ge Zhang, Jian Yang, Tianhao Peng, Zili Wang, Minghao Liu, Junran Peng, Zhaoxiang Zhang, and Jiaheng Liu. Omnivideobench: Towards audio-visual understanding evaluation for omni mllms, 2025a. https://arxiv.org/abs/2510.10689.

Cheng Li, Ziang Leng, Chenxi Yan, Junyi Shen, Hao Wang, Weishi Mi, Yaying Fei, Xiaoyang Feng, Song Yan, HaoSheng Wang, Linkang Zhan, Yaokai Jia, Pingyu Wu, and Haozhen Sun. Chatharuhi: Reviving anime character in reality via large language model. *CoRR*, abs/2308.09597, 2023a. doi: 10.48550/ARXIV.2308.09597. https://doi.org/10.48550/arXiv.2308.09597.

Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. CAMEL: Communicative agents for "mind" exploration of large language model society. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023b.

Hao Li, Chenghao Yang, An Zhang, Yang Deng, Xiang Wang, and Tat-Seng Chua. Hello again! llm-powered personalized agent for long-term dialogue, 2025b. https://arxiv.org/abs/2406.05925.

Haoyang Li, Yiming Li, Anxin Tian, Tianhao Tang, Zhanchao Xu, Xuejia Chen, Nicole Hu, Wei Dong, Qing Li, and Lei Chen. A survey on large language model acceleration based on KV cache management. *Trans. Mach. Learn. Res.*, 2025, 2025c. https://openreview.net/forum?id=z3JZzu9EA3.

Jiaang Li, Quan Wang, Zhongnan Wang, Yongdong Zhang, and Zhendong Mao. ELDER: enhancing lifelong model editing with mixture-of-lora. In Toby Walsh, Julie Shah, and Zico Kolter, editors, *AAAI-25, Sponsored by the Association for the Advancement of Artificial Intelligence, February 25 - March 4, 2025, Philadelphia, PA, USA*, pages 24440–24448. AAAI Press, 2025d. doi: 10.1609/AAAI.V39I23.34622. https://doi.org/10.1609/aaai.v39i23.34622.

Kuan Li, Zhongwang Zhang, Huifeng Yin, Liwen Zhang, Litu Ou, Jialong Wu, Wenbiao Yin, Baixuan Li, Zhengwei Tao, Xinyu Wang, Weizhou Shen, Junkai Zhang, Dingchu Zhang, Xixi Wu, Yong Jiang, Ming Yan, Pengjun Xie, Fei Huang, and Jingren Zhou. Websailor: Navigating super-human reasoning for web agent. *CoRR*, abs/2507.02592, 2025e. doi: 10.48550/ARXIV.2507.02592. https://doi.org/10.48550/arXiv.2507.02592.

Rui Li, Zeyu Zhang, Xiaohe Bo, Zihang Tian, Xu Chen, Quanyu Dai, Zhenhua Dong, and Ruiming Tang. Cam: A constructivist view of agentic memory for llm-based reading comprehension, October 2025f.

Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang, Yujia Zhou, Yutao Zhu, Peitian Zhang, and Zhicheng Dou. Search-o1: Agentic search-enhanced large reasoning models. *CoRR*, abs/2501.05366, 2025g. doi: 10.48550/ARXIV.2501.05366. https://doi.org/10.48550/arXiv.2501.05366.

Xiaoxi Li, Wenxiang Jiao, Jiarui Jin, Guanting Dong, Jiajie Jin, Yinuo Wang, Hao Wang, Yutao Zhu, Ji-Rong Wen, Yuan Lu, and Zhicheng Dou. DeepAgent: A General Reasoning Agent with Scalable Toolsets, October 2025h.

Xiaoxi Li, Jiajie Jin, Guanting Dong, Hongjin Qian, Yutao Zhu, Yongkang Wu, Ji-Rong Wen, and Zhicheng Dou. Webthinker: Empowering large reasoning models with deep research capability. *CoRR*, abs/2504.21776, 2025i. doi: 10.48550/ARXIV.2504.21776. https://doi.org/10.48550/arXiv.2504.21776.

Xiaoxi Li, Jiajie Jin, Yujia Zhou, Yuyao Zhang, Peitian Zhang, Yutao Zhu, and Zhicheng Dou. From matching to generation: A survey on generative information retrieval. *ACM Trans. Inf. Syst.*, 43(3):83:1–83:62, 2025j. doi: 10.1145/3722552. https://doi.org/10.1145/3722552.

Yucheng Li, Bo Dong, Frank Guerin, and Chenghua Lin. Compressing context to enhance inference efficiency of large language models. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6342–6353, Singapore, December 2023c. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.391.

Yuhong Li, Yingbing Huang, Bowen Yang, Bharat Venkitesh, Acyr Locatelli, Hanchen Ye, Tianle Cai, Patrick Lewis, and Deming Chen. Snapkv: LLM knows what you are looking for before generation. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang, editors, *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024b. http://papers.nips.cc/paper_files/paper/2024/hash/28ab418242603e0f7323e54185d19bde-Abstract-Conference.html.

Yunxuan Li, Yibing Du, Jiageng Zhang, Le Hou, Peter Grabowski, Yeqing Li, and Eugene Ie. Improving multi-agent debate with sparse communication topology. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024*, pages 7281–7294. Association for Computational Linguistics, 2024c. doi: 10.18653/V1/2024.FINDINGS-EMNLP.427. https://doi.org/10.18653/v1/2024.findings-emnlp.427.

Zaijing Li, Yuquan Xie, Rui Shao, Gongwei Chen, Dongmei Jiang, and Liqiang Nie. Optimus-1: Hybrid multimodal memory empowered agents excel in long-horizon tasks, 2024d. https://arxiv.org/abs/2408.03615.

Zhiyu Li, Shichao Song, Hanyu Wang, Simin Niu, Ding Chen, Jiawei Yang, Chenyang Xi, Huayi Lai, Jihao Zhao, Yezhaohui Wang, Junpeng Ren, Zehao Lin, Jiahao Huo, Tianyi Chen, Kai Chen, Kehang Li, Zhiqiang Yin, Qingchen Yu, Bo Tang, Hongkang Yang, Zhi-Qin John Xu, and Feiyu Xiong. Memos: An operating system for memory-augmented generation (mag) in large language models, 2025k. https://arxiv.org/abs/2505.22101.

Zijian Li, Xin Guan, Bo Zhang, Shen Huang, Houquan Zhou, Shaopeng Lai, Ming Yan, Yong Jiang, Pengjun Xie, Fei Huang, Jun Zhang, and Jingren Zhou. Webweaver: Structuring web-scale evidence with dynamic outlines for open-ended deep research, October 2025l.

Xuechen Liang, Meiling Tao, Yinghui Xia, Jianhui Wang, Kun Li, Yijin Wang, Yangfan He, Jingsong Yang, Tianyu Shi, Yuantao Wang, Miao Zhang, and Xueqian Wang. SAGE: self-evolving agents with reflective and memory-augmented abilities. *Neurocomputing*, 647:130470, 2025. doi: 10.1016/J.NEUCOM.2025.130470. https://doi.org/10.1016/j.neucom.2025.130470.

Huanxuan Liao, Wen Hu, Yao Xu, Shizhu He, Jun Zhao, and Kang Liu. Beyond Hard and Soft: Hybrid Context Compression for Balancing Local and Global Information Retention. *CoRR*, abs/2505.15774, 2025a. doi: 10.48550/ ARXIV.2505.15774.

Yue Liao, Pengfei Zhou, Siyuan Huang, Donglin Yang, Shengcong Chen, Yuxin Jiang, Yue Hu, Jingbin Cai, Si Liu, Jianlan Luo, et al. Genie envisioner: A unified world foundation platform for robotic manipulation. *arXiv preprint arXiv:2508.05635*, 2025b.

Opher Lieber, Barak Lenz, Hofit Bata, Gal Cohen, Jhonathan Osin, Itay Dalmedigos, Erez Safahi, Shaked Meirom, Yonatan Belinkov, Shai Shalev-Shwartz, Omri Abend, Raz Alon, Tomer Asida, Amir Bergman, Roman Glozman, Michael Gokhman, Avashalom Manevich, Nir Ratner, Noam Rozen, Erez Shwartz, Mor Zusman, and Yoav Shoham. Jamba: A hybrid transformer-mamba language model, 2024. https://arxiv.org/abs/2403.19887.

Thomas Limbacher and Robert Legenstein. H-mem: Harnessing synaptic plasticity with hebbian memory networks. In *Advances in Neural Information Processing Systems*, volume 33, pages 21627–21637. Curran Associates, Inc., 2020.

Jessy Lin, Luke Zettlemoyer, Gargi Ghosh, Wen-Tau Yih, Aram Markosyan, Vincent-Pierre Berges, and Barlas Oğuz. Continual learning via sparse memory finetuning, 2025. https://arxiv.org/abs/2510.15103.

Guangyi Liu, Pengxiang Zhao, Liang Liu, Zhiming Chen, Yuxiang Chai, Shuai Ren, Hao Wang, Shibo He, and Wenchao Meng. Learnact: Few-shot mobile GUI agent with a unified demonstration benchmark. *CoRR*, abs/2504.13805, 2025a. doi: 10.48550/ARXIV.2504.13805. https://doi.org/10.48550/arXiv.2504.13805.

Jiahao Liu, Shengkang Gu, Dongsheng Li, Guangping Zhang, Mingzhe Han, Hansu Gu, Peng Zhang, Tun Lu, Li Shang, and Ning Gu. Agentcf++: Memory-enhanced llm-based agents for popularity-aware cross-domain recommendations, 2025b. https://arxiv.org/abs/2502.13843.

Jiale Liu, Yifan Zeng, Malte Højmark-Bertelsen, Marie Normann Gadeberg, Huazheng Wang, and Qingyun Wu. Memory-augmented agent training for business document understanding. *CoRR*, abs/2412.15274, 2024. doi: 10.48550/ARXIV.2412.15274. https://doi.org/10.48550/arXiv.2412.15274.

Jun Liu, Zhenglun Kong, Changdi Yang, Fan Yang, Tianqi Li, Peiyan Dong, Joannah Nanjekye, Hao Tang, Geng Yuan, Wei Niu, Wenbin Zhang, Pu Zhao, Xue Lin, Dong Huang, and Yanzhi Wang. Rcr-router: Efficient role-aware context routing for multi-agent LLM systems with structured memory. *CoRR*, abs/2508.04903, 2025c. doi: 10.48550/ARXIV.2508.04903. https://doi.org/10.48550/arXiv.2508.04903.

Junming Liu, Yifei Sun, Weihua Cheng, Haodong Lei, Yirong Chen, Licheng Wen, Xuemeng Yang, Daocheng Fu, Pinlong Cai, Nianchen Deng, Yi Yu, Shuyue Hu, Botian Shi, and Ding Wang. Memverse: Multimodal memory for lifelong learning agents, 2025d. https://arxiv.org/abs/2512.03627.

Kunhao Liu, Wenbo Hu, Jiale Xu, Ying Shan, and Shijian Lu. Rolling forcing: Autoregressive long video diffusion in real time. *arXiv preprint arXiv:2509.25161*, 2025e.

Lei Liu, Xiaoyan Yang, Yue Shen, Binbin Hu, Zhiqiang Zhang, Jinjie Gu, and Guannan Zhang. Think-in-memory: Recalling and post-thinking enable llms with long-term memory. *CoRR*, abs/2311.08719, 2023a. doi: 10.48550/ ARXIV.2311.08719. https://doi.org/10.48550/arXiv.2311.08719.

Yanming Liu, Xinyue Peng, Jiannan Cao, Shi Bo, Yuwei Zhang, Xuhong Zhang, Sheng Cheng, Xun Wang, Jianwei Yin, and Tianyu Du. Tool-Planner: Task planning with clusters across multiple tools. In *The Thirteenth International Conference on Learning Representations*, 2025f.

Yixin Liu, Guibin Zhang, Kun Wang, Shiyuan Li, and Shirui Pan. Graph-augmented large language model agents: Current progress and future prospects. *arXiv preprint arXiv:2507.21407*, 2025g.

Zeting Liu, Zida Yang, Zeyu Zhang, and Hao Tang. Evovla: Self-evolving vision-language-action model, 2025h. https://arxiv.org/abs/2511.16166.

Zichang Liu, Aditya Desai, Fangshuo Liao, Weitao Wang, Victor Xie, Zhaozhuo Xu, Anastasios Kyrillidis, and Anshumali Shrivastava. Scissorhands: Exploiting the persistence of importance hypothesis for LLM KV cache compression at test time. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023b. http://papers.nips.cc/paper_files/paper/2023/hash/a452a7c6c463e4ae8fbdc614c6e983e6-Abstract-Conference.html.

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

Junfeng Lu and Yueyan Li. Dynamic affective memory management for personalized llm agents, 2025. https://arxiv.org/abs/2510.27418.

Junru Lu, Siyu An, Mingbao Lin, Gabriele Pergola, Yulan He, Di Yin, Xing Sun, and Yunsheng Wu. MemoChat: Tuning LLMs to use memos for consistent long-range open-domain conversation. *arXiv preprint arXiv:2308.08239*, 2023.

Miao Lu, Weiwei Sun, Weihua Du, Zhan Ling, Xuesong Yao, Kang Liu, and Jiecao Chen. Scaling llm multi-turn rl with end-to-end summarization-based context management, 2025a. https://arxiv.org/abs/2510.06727.

Miao Lu, Weiwei Sun, Weihua Du, Zhan Ling, Xuesong Yao, Kang Liu, and Jiecao Chen. Scaling llm multi-turn rl with end-to-end summarization-based context management, October 2025b.

Pengqian Lu, Jie Lu, Anjin Liu, and Guangquan Zhang. Spad: Seven-source token probability attribution with syntactic aggregation for detecting hallucinations in rag. *arXiv preprint arXiv:2512.07515*, 2025c.

Elias Lumer, Anmol Gulati, Vamse Kumar Subbiah, Pradeep Honaganahalli Basavaraju, and James A Burke. Memtool: Optimizing short-term memory management for dynamic tool calling in llm agent multi-turn conversations. *arXiv preprint arXiv:2507.21428*, 2025.

Junyu Luo, Weizhi Zhang, Ye Yuan, Yusheng Zhao, Junwei Yang, Yiyang Gu, Bohan Wu, Binqi Chen, Ziyue Qiao, Qingqing Long, Rongcheng Tu, Xiao Luo, Wei Ju, Zhiping Xiao, Yifan Wang, Meng Xiao, Chenwu Liu, Jingyang Yuan, Shichang Zhang, Yiqiao Jin, Fan Zhang, Xian Wu, Hanqing Zhao, Dacheng Tao, Philip S. Yu, and Ming Zhang. Large language model agent: A survey on methodology, applications and challenges, 2025. https://arxiv.org/abs/2503.21460.

Weiyao Luo, Suncong Zheng, Heming Xia, Weikang Wang, Yan Lei, Tianyu Liu, Shuang Chen, and Zhifang Sui. Taking a deep breath: Enhancing language modeling of large language models with sentinel tokens. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024*, pages 4034–4040. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.FINDINGS-EMNLP.233. https://doi.org/10.18653/v1/2024.findings-emnlp.233.

Huan Ma, Changqing Zhang, Yatao Bian, Lemao Liu, Zhirui Zhang, Peilin Zhao, Shu Zhang, Huazhu Fu, Qinghua Hu, and Bingzhe Wu. Fairness-guided few-shot prompting for large language models, 2023a. https://arxiv.org/abs/2303.13217.

Jun-Yu Ma, Zhen-Hua Ling, Ningyu Zhang, and Jia-Chen Gu. Neighboring perturbations of knowledge editing on large language models. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024. https://openreview.net/forum?id=K9NTPRvVRI.

Xinbei Ma, Yeyun Gong, Pengcheng He, Hai Zhao, and Nan Duan. Query rewriting in retrieval-augmented large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5303–5315, 2023b.

Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-Refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, pages 46534–46594, 2023.

Adyasha Maharana, Dong-Ho Lee, Sergey Tulyakov, Mohit Bansal, Francesco Barbieri, and Yuwei Fang. Evaluating very long-term conversational memory of LLM agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13851–13870, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.747. https://aclanthology.org/2024.acl-long.747/.

Shengyu Mao, Xiaohan Wang, Mengru Wang, Yong Jiang, Pengjun Xie, Fei Huang, and Ningyu Zhang. Editing personality for large language models. In Derek F. Wong, Zhongyu Wei, and Muyun Yang, editors, *Natural Language Processing and Chinese Computing - 13th National CCF Conference, NLPCC 2024, Hangzhou, China, November 1-3, 2024, Proceedings, Part II*, volume 15360 of *Lecture Notes in Computer Science*, pages 241–254. Springer, 2024. doi: 10.1007/978-981-97-9434-8\_19. https://doi.org/10.1007/978-981-97-9434-8_19.

Samuele Marro, Emanuele La Malfa, Jesse Wright, Guohao Li, Nigel Shadbolt, Michael Wooldridge, and Philip Torr. A scalable communication protocol for networks of large language models, 2024. https://arxiv.org/abs/2410.11905.

Andrea Matarazzo and Riccardo Torlone. A survey on large language models with some insights on their capabilities and limitations, 2025. https://arxiv.org/abs/2501.04040.

Marcelo G. Mattar and Nathaniel D. Daw. Prioritized memory access explains planning and hippocampal replay. *Nature Neuroscience*, 21(11):1609–1617, November 2018. ISSN 1546-1726. doi: 10.1038/s41593-018-0232-z.

James L. McClelland, Bruce L. McNaughton, and Randall C. O'Reilly. Why there are complementary learning systems in the hippocampus and neocortex: Insights from the successes and failures of connectionist models of learning and memory. *Psychological Review*, 102(3):419–457, July 1995. ISSN 0033-295X. doi: 10.1037/0033-295X.102.3.419.

Michael McCloskey and Neal J. Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. volume 24 of *Psychology of Learning and Motivation*, pages 109–165. Academic Press, 1989. doi: https://doi.org/10.1016/S0079-7421(08)60536-8. https://www.sciencedirect.com/science/article/pii/S0079742108605368.

Lingrui Mei, Jiayu Yao, Yuyao Ge, Yiwei Wang, Baolong Bi, Yujun Cai, Jiazhi Liu, Mingyu Li, Zhong-Zhi Li, Duzhen Zhang, Chenlin Zhou, Jiayi Mao, Tianze Xia, Jiafeng Guo, and Shenghua Liu. A Survey of Context Engineering for Large Language Models, July 2025.

Memary. GitHub - kingjulio8238/Memary: The Open Source Memory Layer For Autonomous Agents — github.com. https://github.com/kingjulio8238/Memary, 2025. [Accessed 14-12-2025].

Memobase. GitHub - memodb-io/memobase: User Profile-Based Long-Term Memory for AI Chatbot Applications. https://github.com/memodb-io/memobase, 2025. [Accessed 12-12-2025].

Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in GPT. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. http://papers.nips.cc/paper_files/paper/2022/hash/6f1d43d5a82a37e89b0665b33bf3a182-Abstract-Conference.html.

Kevin Meng, Arnab Sen Sharma, Alex J. Andonian, Yonatan Belinkov, and David Bau. Mass-editing memory in a transformer. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. https://openreview.net/forum?id=MkbcAHIYgyS.

Lina Mezghani, Sainbayar Sukhbaatar, Thibaut Lavril, Oleksandr Maksymets, Dhruv Batra, Piotr Bojanowski, and Karteek Alahari. Memory-augmented reinforcement learning for image-goal navigation. In *IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2022, Kyoto, Japan, October 23-27, 2022*, pages 3316–3323. IEEE, 2022. doi: 10.1109/IROS47612.2022.9981090. https://doi.org/10.1109/IROS47612.2022.9981090.

Grégoire Mialon, Clémentine Fourrier, Thomas Wolf, Yann LeCun, and Thomas Scialom. GAIA: a benchmark for general ai assistants. In *The Twelfth International Conference on Learning Representations*, 2023.

Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. Large language models: A survey, 2025. https://arxiv.org/abs/2402.06196.

MineContext. GitHub - volcengine/MineContext: MineContext is your proactive context-aware AI partner (Context-Engineering+ChatGPT Pulse) — github.com. https://github.com/volcengine/MineContext, 2025. [Accessed 14-12-2025].

Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. Fast model editing at scale. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. https://openreview.net/forum?id=0DcZxeWfOPt.

Atsuyuki Miyai, Zaiying Zhao, Kazuki Egashira, Atsuki Sato, Tatsumi Sunada, Shota Onohara, Hiromasa Yamanishi, Mashiro Toyooka, Kunato Nishina, Ryoma Maeda, et al. Webchorearena: Evaluating web browsing agents on realistic tedious web tasks. *arXiv preprint arXiv:2506.01952*, 2025.

Ali Modarressi, Ayyoob Imani, Mohsen Fayyaz, and Hinrich Schütze. RET-LLM: Towards a general read-write memory for large language models. *arXiv preprint arXiv:2305.14322*, 2023.

Sajad Mousavi, Ricardo Luna Gutiérrez, Desik Rengarajan, Vineet Gundecha, Ashwin Ramesh Babu, Avisek Naug, Antonio Guillen, and Soumyendu Sarkar. N-critics: Self-refinement of large language models with ensemble of critics, 2023. https://arxiv.org/abs/2310.18679.

Jesse Mu, Xiang Li, and Noah D. Goodman. Learning to compress prompts with gist tokens. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. http://papers.nips.cc/paper_files/paper/2023/hash/3d77c6dcc7f143aa2154e7f4d5e22d68-Abstract-Conference.html.

Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. Orca: Progressive learning from complex explanation traces of gpt-4, 2023. https://arxiv.org/abs/2306.02707.

Hyungho Na, Yunkyeong Seo, and Il-Chul Moon. Efficient episodic memory utilization of cooperative multi-agent reinforcement learning. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. https://openreview.net/forum?id=LjivA1SLZ6.

Jiayan Nan, Wenquan Ma, Wenlong Wu, and Yize Chen. Nemori: Self-organizing agent memory inspired by cognitive science. *CoRR*, abs/2508.03341, 2025. doi: 10.48550/ARXIV.2508.03341. https://doi.org/10.48550/arXiv.2508.03341.

Thang Nguyen, Peter Chin, and Yu-Wing Tai. MA-RAG: multi-agent retrieval-augmented generation via collaborative chain-of-thought reasoning. *CoRR*, abs/2505.20096, 2025. doi: 10.48550/ARXIV.2505.20096. https://doi.org/10.48550/arXiv.2505.20096.

Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang, Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin, Ji-Rong Wen, and Chongxuan Li. Large language diffusion models, 2025. https://arxiv.org/abs/2502.09992.

Cheng Niu, Yuanhao Wu, Juno Zhu, Siliang Xu, Kashun Shum, Randy Zhong, Juntong Song, and Tong Zhang. Ragtruth: A hallucination corpus for developing trustworthy retrieval-augmented language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10862–10878, 2024.

OpenAI. Memory and new controls for chatgpt, 2024. https://openai.com/index/memory-and-new-controls-for-chatgpt/.

Siru Ouyang, Jun Yan, I-Hung Hsu, Yanfei Chen, Ke Jiang, Zifeng Wang, Rujun Han, Long T. Le, Samira Daruki, Xiangru Tang, Vishy Tirumalashetty, George Lee, Mahsan Rofouei, Hangfei Lin, Jiawei Han, Chen-Yu Lee, and Tomas Pfister. Reasoningbank: Scaling agent self-evolving with reasoning memory, 2025. https://arxiv.org/abs/2509.25140.

Charles Packer, Vivian Fang, Shishir G. Patil, Kevin Lin, Sarah Wooders, and Joseph E. Gonzalez. MemGPT: Towards LLMs as Operating Systems. *CoRR*, abs/2310.08560, 2023a. doi: 10.48550/ARXIV.2310.08560.

Charles Packer, Vivian Fang, ShishirG Patil, Kevin Lin, Sarah Wooders, and JosephE Gonzalez. Memgpt: Towards llms as operating systems. *arXiv preprint arXiv:2310.08560*, 2023b.

Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The pagerank citation ranking : Bringing order to the web. In *The Web Conference*, 1999. https://api.semanticscholar.org/CorpusID:1508503.

Yiyuan Pan, Yunzhe Xu, Zhe Liu, and Hesheng Wang. Planning from imagination: Episodic simulation and episodic memory for vision-and-language navigation, 2024. https://arxiv.org/abs/2412.01857.

Zhuoshi Pan, Qianhui Wu, Huiqiang Jiang, Xufang Luo, Hao Cheng, Dongsheng Li, Yuqing Yang, Chin-Yew Lin, H. Vicky Zhao, Lili Qiu, and Jianfeng Gao. Secom: On memory construction and retrieval for personalized conversational agents. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net, 2025. https://openreview.net/forum?id=xKDZAW0He3.

Joon Sung Park, Joseph O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the 36th annual acm symposium on user interface software and technology*, pages 1–22, 2023.

Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. Gorilla: Large Language Model Connected with Massive APIs. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang, editors, *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024. http://papers.nips.cc/paper_files/paper/2024/hash/e4c61f578ff07830f5c37378dd3ecb0d-Abstract-Conference.html.

Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, Huanqi Cao, Xin Cheng, Michael Chung, Matteo Grella, Kranthi Kiran GV, Xuzheng He, Haowen Hou, Jiaju Lin, Przemyslaw Kazienko, Jan Kocon, Jiaming Kong, Bartlomiej Koptyra, Hayden Lau, Krishna Sri Ipsit Mantri, Ferdinand Mom, Atsushi Saito, Guangyu Song, Xiangru Tang, Bolun Wang, Johan S. Wind, Stanislaw Wozniak, Ruichong Zhang, Zhenyuan Zhang, Qihang Zhao, Peng Zhou, Qinghua Zhou, Jian Zhu, and Rui-Jie Zhu. Rwkv: Reinventing rnns for the transformer era, 2023. https://arxiv.org/abs/2305.13048.

Boci Peng, Yun Zhu, Yongchao Liu, Xiaohe Bo, Haizhou Shi, Chuntao Hong, Yan Zhang, and Siliang Tang. Graph retrieval-augmented generation: A survey, 2024. https://arxiv.org/abs/2408.08921.

Jayr Alencar Pereira, Robson do Nascimento Fidalgo, Roberto A. Lotufo, and Rodrigo Nogueira. Visconde: Multi-document QA with GPT-3 and neural reranking. In Jaap Kamps, Lorraine Goeuriot, Fabio Crestani, Maria Maistro, Hideo Joho, Brian Davis, Cathal Gurrin, Udo Kruschwitz, and Annalina Caputo, editors, *Advances in Information Retrieval - 45th European Conference on Information Retrieval, ECIR 2023, Dublin, Ireland, April 2-6, 2023, Proceedings, Part II*, volume 13981 of *Lecture Notes in Computer Science*, pages 534–543. Springer, 2023. doi: 10.1007/978-3-031-28238-6\_44. https://doi.org/10.1007/978-3-031-28238-6_44.

Jinghua Piao, Yuwei Yan, Jun Zhang, Nian Li, Junbo Yan, Xiaochong Lan, Zhihong Lu, Zhiheng Zheng, Jing Yi Wang, Di Zhou, Chen Gao, Fengli Xu, Fang Zhang, Ke Rong, Jun Su, and Yong Li. Agentsociety: Large-scale simulation of llm-driven generative agents advances understanding of human behaviors and society, February 2025.

Ryan Po, Yotam Nitzan, Richard Zhang, Berlin Chen, Tri Dao, Eli Shechtman, Gordon Wetzstein, and Xun Huang. Long-context state-space video world models. *arXiv preprint arXiv:2505.20171*, 2025.

Hadi Pouransari, David Grangier, C Thomas, Michael Kirchhof, and Oncel Tuzel. Pretraining with hierarchical memories: separating long-tail and common knowledge, 2025. https://arxiv.org/abs/2510.02375.

Shrimai Prabhumoye, Rafal Kocielnik, Mohammad Shoeybi, Anima Anandkumar, and Bryan Catanzaro. Few-shot instruction prompts for pretrained language models to detect social biases, 2022. https://arxiv.org/abs/2112.07868.

Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. ChatDev: Communicative agents for software development. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15174–15186, 2024.

Cheng Qian, Chi Han, Yi Ren Fung, Yujia Qin, Zhiyuan Liu, and Heng Ji. CREATOR: Tool Creation for Disentangling Abstract and Concrete Reasoning of Large Language Models. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 6922–6939. Association for Computational Linguistics, 2023. doi: 10.18653/V1/2023.FINDINGS-EMNLP.462. https://doi.org/10.18653/v1/2023.findings-emnlp.462.

Hongjin Qian, Zheng Liu, Peitian Zhang, Kelong Mao, Defu Lian, Zhicheng Dou, and Tiejun Huang. Memorag: Boosting long context processing with global memory-enhanced retrieval augmentation. In Guodong Long, Michale Blumestein, Yi Chang, Liane Lewin-Eytan, Zi Helen Huang, and Elad Yom-Tov, editors, *Proceedings of the ACM on Web Conference 2025, WWW 2025, Sydney, NSW, Australia, 28 April 2025- 2 May 2025*, pages 2366–2377. ACM, 2025. doi: 10.1145/3696410.3714805. https://doi.org/10.1145/3696410.3714805.

Tianrui Qin, Qianben Chen, Sinuo Wang, He Xing, King Zhu, He Zhu, Dingfeng Shi, Xinxin Liu, Ge Zhang, Jiaheng Liu, Yuchen Eleanor Jiang, Xitong Gao, and Wangchunshu Zhou. Flash-searcher: Fast and effective web agents via dag-based parallel execution, 2025. https://arxiv.org/abs/2509.25301.

Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. ToolLLM: Facilitating large language models to master 16000+ real-world apis. In *The Twelfth International Conference on Learning Representations*, 2024a.

Zhen Qin, Weigao Sun, Dong Li, Xuyang Shen, Weixuan Sun, and Yiran Zhong. Various lengths, constant speed: Efficient language modeling with lightning attention. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024b. https://openreview.net/forum?id=Lwm6TiUP4X.

Zhen Qin, Weigao Sun, Dong Li, Xuyang Shen, Weixuan Sun, and Yiran Zhong. Lightning attention-2: A free lunch for handling unlimited sequence lengths in large language models. *CoRR*, abs/2401.04658, 2024c. doi: 10.48550/ARXIV.2401.04658. https://doi.org/10.48550/arXiv.2401.04658.

Jiahao Qiu, Xinzhe Juan, Yimin Wang, Ling Yang, Xuan Qi, Tongcheng Zhang, Jiacheng Guo, Yifu Lu, Zixin Yao, Hongru Wang, Shilong Liu, Xun Jiang, Liu Leqi, and Mengdi Wang. Agentdistill: Training-free agent distillation with generalizable mcp boxes, 2025a. https://arxiv.org/abs/2506.14728.

Jiahao Qiu, Xuan Qi, Hongru Wang, Xinzhe Juan, Yimin Wang, Zelin Zhao, Jiayi Geng, Jiacheng Guo, Peihang Li, Jingzhe Shi, Shilong Liu, and Mengdi Wang. Alita-g: Self-evolving generative agent for agent generation, October 2025b.

Jiahao Qiu, Xuan Qi, Tongcheng Zhang, Xinzhe Juan, Jiacheng Guo, Yifu Lu, Yimin Wang, Zixin Yao, Qihan Ren, Xun Jiang, Xing Zhou, Dongrui Liu, Ling Yang, Yue Wu, Kaixuan Huang, Shilong Liu, Hongru Wang, and Mengdi Wang.

Alita: Generalist agent enabling scalable agentic reasoning with minimal predefinition and maximal self-evolution, 2025c. https://arxiv.org/abs/2505.20286.

Changle Qu, Sunhao Dai, Xiaochi Wei, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, Jun Xu, and Ji-Rong Wen. COLT: Towards Completeness-Oriented Tool Retrieval for Large Language Models. *CoRR*, abs/2405.16089, 2024. doi: 10.48550/ARXIV.2405.16089. https://doi.org/10.48550/arXiv.2405.16089. arXiv: 2405.16089.

Changle Qu, Sunhao Dai, Xiaochi Wei, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, Jun Xu, and Ji-Rong Wen. From exploration to mastery: Enabling llms to master tools via self-driven interactions, February 2025a.

Changle Qu, Sunhao Dai, Xiaochi Wei, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, Jun Xu, and Ji-rong Wen. Tool learning with large language models: a survey. *Frontiers of Computer Science*, 19(8), January 2025b. ISSN 2095-2236. doi: 10.1007/s11704-024-40678-2. http://dx.doi.org/10.1007/s11704-024-40678-2.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR, 2021. http://proceedings.mlr.press/v139/radford21a.html.

Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.

Krishan Rana, Jesse Haviland, Sourav Garg, Jad Abou-Chakra, Ian D. Reid, and Niko Sünderhauf. SayPlan: Grounding Large Language Models using 3D Scene Graphs for Scalable Robot Task Planning. In Jie Tan, Marc Toussaint, and Kourosh Darvish, editors, *Conference on Robot Learning, CoRL 2023, 6-9 November 2023, Atlanta, GA, USA*, volume 229 of *Proceedings of Machine Learning Research*, pages 23–72. PMLR, 2023.

Preston Rasmussen, Pavlo Paliychuk, Travis Beauvais, Jack Ryan, and Daniel Chalef. Zep: A temporal knowledge graph architecture for agent memory. *CoRR*, abs/2501.13956, 2025. doi: 10.48550/ARXIV.2501.13956. https://doi.org/10.48550/arXiv.2501.13956.

Paul J. Reber. The neural basis of implicit learning and memory: A review of neuropsychological and neuroimaging research. *Neuropsychologia*, 51(10):2026–2042, August 2013. ISSN 1873-3514. doi: 10.1016/j.neuropsychologia.2013.06.019.

Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pages 3980–3990. Association for Computational Linguistics, 2019. doi: 10.18653/V1/D19-1410. https://doi.org/10.18653/v1/D19-1410.

Alireza Rezazadeh, Zichao Li, Ange Lou, Yuying Zhao, Wei Wei, and Yujia Bao. Collaborative memory: Multi-user memory sharing in llm agents with dynamic access control. *arXiv preprint arXiv:2505.18279*, 2025a.

Alireza Rezazadeh, Zichao Li, Ange Lou, Yuying Zhao, Wei Wei, and Yujia Bao. Collaborative memory: Multi-user memory sharing in llm agents with dynamic access control, 2025b. https://arxiv.org/abs/2505.18279.

Alireza Rezazadeh, Zichao Li, Wei Wei, and Yujia Bao. From Isolated Conversations to Hierarchical Schemas: Dynamic Tree Memory Representation for LLMs. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net, 2025c.

Wim J. Riedel and Arjan Blokland. Declarative memory. *Handbook of Experimental Pharmacology*, 228:215–236, 2015. ISSN 0171-2004. doi: 10.1007/978-3-319-16522-6_7.

Stephen E. Robertson and Hugo Zaragoza. The probabilistic relevance framework: BM25 and beyond. *Found. Trends Inf. Retr.*, 3(4):333–389, 2009. doi: 10.1561/1500000019. https://doi.org/10.1561/1500000019.

Dongyu Ru, Lin Qiu, Xiangkun Hu, Tianhang Zhang, Peng Shi, Shuaichen Chang, Cheng Jiayang, Cunxiang Wang, Shichao Sun, Huanyu Li, et al. Ragchecker: A fine-grained framework for diagnosing retrieval-augmented generation. *Advances in Neural Information Processing Systems*, 37:21999–22027, 2024.

Jingqing Ruan, Yihong Chen, Bin Zhang, Zhiwei Xu, Tianpeng Bao, Guoqing Du, Shiwei Shi, Hangyu Mao, Ziyue Li, Xingyu Zeng, and Rui Zhao. Tptu: Large language model-based ai agents for task planning and tool usage, 2023. https://arxiv.org/abs/2308.03427.

Rana Salama, Jason Cai, Michelle Yuan, Anna Currey, Monica Sunkara, Yi Zhang, and Yassine Benajiba. Meminsight: Autonomous memory augmentation for llm agents, 2025. https://arxiv.org/abs/2503.21760.

Jitao Sang, Jinlin Xiao, Jiarun Han, Jilin Chen, Xiaoyi Chen, Shuyu Wei, Yongjie Sun, and Yuhang Wang. Beyond pipelines: A survey of the paradigm shift toward model-native agentic ai, 2025. https://arxiv.org/abs/2510.16720.

Parth Sarthi, Salman Abdullah, Aditi Tuli, Shubh Khanna, Anna Goldie, and Christopher D. Manning. Raptor: Recursive abstractive processing for tree-organized retrieval. *ArXiv*, abs/2401.18059, 2024.

Daniel L. Schacter and Donna Rose Addis. Constructive memory: The ghosts of past and future. *Nature*, 445(7123): 27, January 2007. ISSN 1476-4687. doi: 10.1038/445027a.

Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. ToolFormer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:68539–68551, 2023.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

Carol A. Seger and Brian J. Spiering. A critical review of habit learning and the Basal Ganglia. *Frontiers in Systems Neuroscience*, 5:66, 2011. ISSN 1662-5137. doi: 10.3389/fnsys.2011.00066.

Jay Shah, Ganesh Bikshandi, Ying Zhang, Vijay Thakkar, Pradeep Ramani, and Tri Dao. Flashattention-3: Fast and accurate attention with asynchrony and low-precision. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang, editors, *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024. http://papers.nips.cc/paper_files/paper/2024/hash/7ede97c3e082c6df10a8d6103a2eebd2-Abstract-Conference.html.

Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. Character-llm: A trainable agent for role-playing. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 13153–13187. Association for Computational Linguistics, 2023. doi: 10.18653/V1/2023.EMNLP-MAIN.814. https://doi.org/10.18653/v1/2023.emnlp-main.814.

Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.

Junxiao Shen, John J. Dudley, and Per Ola Kristensson. Encode-Store-Retrieve: Augmenting Human Memory through Language-Encoded Egocentric Perception. In Ulrich Eck, Misha Sra, Jeanine K. Stefanucci, Maki Sugimoto, Markus Tatzgern, and Ian Williams, editors, *IEEE International Symposium on Mixed and Augmented Reality, ISMAR 2024, Bellevue, WA, USA, October 21-25, 2024*, pages 923–931. IEEE, 2024. doi: 10.1109/ISMAR62088.2024.00108.

Hao Shi, Bin Xie, Yingfei Liu, Lin Sun, Fengrong Liu, Tiancai Wang, Erjin Zhou, Haoqiang Fan, Xiangyu Zhang, and Gao Huang. Memoryvla: Perceptual-cognitive memory in vision-language-action models for robotic manipulation. *CoRR*, abs/2508.19236, 2025a. doi: 10.48550/ARXIV.2508.19236. https://doi.org/10.48550/arXiv.2508.19236.

Yaorui Shi, Yuxin Chen, Siyuan Wang, Sihang Li, Hengxing Cai, Qi Gu, Xiang Wang, and An Zhang. Look back to reason forward: Revisitable memory for long-context llm agents, 2025b. https://arxiv.org/abs/2509.23040.

Zhengliang Shi, Yuhan Wang, Lingyong Yan, Pengjie Ren, Shuaiqiang Wang, Dawei Yin, and Zhaochun Ren. Retrieval models aren't tool-savvy: Benchmarking tool retrieval for large language models, May 2025c.

Zitong Shi, Guancheng Wan, Wenke Huang, Guibin Zhang, Jiawei Shao, Mang Ye, and Carl Yang. Privacy-enhancing paradigms within federated multi-agent systems. *arXiv preprint arXiv:2503.08175*, 2025d.

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, 2023a.

Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: language agents with verbal reinforcement learning. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023b. http://papers.nips.cc/paper_files/paper/2023/hash/1b44b878bb782e6954cd888628510e90-Abstract-Conference.html.

Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew J. Hausknecht. Alfworld: Aligning text and embodied environments for interactive learning. In *9th International Conference on Learning Representations*, 2021.

Aditi Singh, Abul Ehtesham, Saket Kumar, and Tala Talaei Khoei. Agentic retrieval-augmented generation: A survey on agentic rag, 2025. https://arxiv.org/abs/2501.09136.

Anton Sinitsin, Vsevolod Plokhotnyuk, Dmitry V. Pyrkin, Sergei Popov, and Artem Babenko. Editable neural networks. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. https://openreview.net/forum?id=HJedXaEtvS.

Chan Hee Song, Brian M. Sadler, Jiaman Wu, Wei-Lun Chao, Clayton Washington, and Yu Su. LLM-Planner: Few-Shot Grounded Planning for Embodied Agents with Large Language Models. In *IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023*, pages 2986–2997. IEEE, 2023. doi: 10.1109/ICCV51070.2023.00280.

Enxin Song, Wenhao Chai, Guanhong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Haozhe Chi, Xun Guo, Tian Ye, Yanting Zhang, Yan Lu, Jenq-Neng Hwang, and Gaoang Wang. Moviechat: From dense token to sparse memory for long video understanding. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024*, pages 18221–18232. IEEE, 2024. doi: 10.1109/CVPR52733.2024.01725. https://doi.org/10.1109/CVPR52733.2024.01725.

Huatong Song, Jinhao Jiang, Yingqian Min, Jie Chen, Zhipeng Chen, Wayne Xin Zhao, Lei Fang, and Ji-Rong Wen. R1-searcher: Incentivizing the search capability in llms via reinforcement learning. *CoRR*, abs/2503.05592, 2025a. doi: 10.48550/ARXIV.2503.05592. https://doi.org/10.48550/arXiv.2503.05592.

Xinshuai Song, Weixing Chen, Yang Liu, Weikai Chen, Guanbin Li, and Liang Lin. Towards long-horizon vision-language navigation: Platform, benchmark and method. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2025, Nashville, TN, USA, June 11-15, 2025*, pages 12078–12088. Computer Vision Foundation / IEEE, 2025b. doi: 10.1109/CVPR52734.2025.01128. https://openaccess.thecvf.com/content/CVPR2025/html/Song_Towards_Long-Horizon_Vision-Language_Navigation_Platform_Benchmark_and_Method_CVPR_2025_paper.html.

KAREN SPARCK JONES. A STATISTICAL INTERPRETATION OF TERM SPECIFICITY AND ITS APPLICATION IN RETRIEVAL. *Journal of Documentation*, 28(1):11–21, January 1972. ISSN 0022-0418. doi: 10.1108/eb026526.

Larry R. Squire. Memory systems of the brain: A brief history and current perspective. *Neurobiology of Learning and Memory*, 82(3):171–177, November 2004. ISSN 1074-7427. doi: 10.1016/j.nlm.2004.06.005.

Liangcai Su, Zhen Zhang, Guangyu Li, Zhuo Chen, Chenxi Wang, Maojia Song, Xinyu Wang, Kuan Li, Jialong Wu, Xuanzhong Chen, Zile Qiao, Zhongwang Zhang, Huifeng Yin, Shihao Cai, Runnan Fang, Zhengwei Tao, Wenbiao Yin, Chenxiong Qian, Yong Jiang, Pengjun Xie, Fei Huang, and Jingren Zhou. Scaling agents via continual pre-training. *CoRR*, abs/2509.13310, 2025. doi: 10.48550/ARXIV.2509.13310. https://doi.org/10.48550/arXiv.2509.13310.

Haoran Sun and Shaoning Zeng. Hierarchical memory for high-efficiency long-term reasoning in llm agents. *ArXiv*, abs/2507.22925, 2025.

Jingwei Sun, Zhixu Du, and Yiran Chen. Knowledge graph tuning: Real-time large language model personalization based on human feedback, 2024. https://arxiv.org/abs/2405.19686.

Weiwei Sun, Miao Lu, Zhan Ling, Kang Liu, Xuesong Yao, Yiming Yang, and Jiecao Chen. Scaling long-horizon LLM agent via context-folding. *CoRR*, abs/2510.11967, 2025a. doi: 10.48550/ARXIV.2510.11967. https://doi.org/10.48550/arXiv.2510.11967.

Zeyi Sun, Ziyu Liu, Yuhang Zang, Yuhang Cao, Xiaoyi Dong, Tong Wu, Dahua Lin, and Jiaqi Wang. SEAgent: Self-evolving computer use agent with autonomous learning from experience. *arXiv preprint arXiv:2508.04700*, 2025b.

Zhongxiang Sun, Zihua Si, Xiaoxue Zang, Kai Zheng, Yang Song, Xiao Zhang, and Jun Xu. Largepig for hallucination-free query generation: Your large language model is secretly a pointer generator. In *Proceedings of the ACM on Web Conference 2025*, WWW '25, page 4766–4779, New York, NY, USA, 2025c. Association for Computing Machinery. ISBN 9798400712746. doi: 10.1145/3696410.3714800. https://doi.org/10.1145/3696410.3714800.

Zhongxiang Sun, Qipeng Wang, Weijie Yu, Xiaoxue Zang, Kai Zheng, Jun Xu, Xiao Zhang, Yang Song, and Han Li.

Rearter: Retrieval-augmented reasoning with trustworthy process rewarding. In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1251–1261, 2025d.

ZhongXiang Sun, Xiaoxue Zang, Kai Zheng, Jun Xu, Xiao Zhang, Weijie Yu, Yang Song, and Han Li. Redeep: Detecting hallucination in retrieval-augmented generation via mechanistic interpretability. In *The Thirteenth International Conference on Learning Representations*, 2025e.

Supermemory. Supermemory — Universal Memory API for AI apps — supermemory.ai. https://supermemory.ai/, 2025. [Accessed 14-12-2025].

David Silver Sutton, Richard S. Welcome to the Era of Experience, April 2025.

Mirac Suzgun, Mert Yuksekgonul, Federico Bianchi, Dan Jurafsky, and James Zou. Dynamic cheatsheet: Test-time learning with adaptive memory, 2025. https://arxiv.org/abs/2504.07952.

Jihoon Tack, Jaehyung Kim, Eric Mitchell, Jinwoo Shin, Yee Whye Teh, and Jonathan Richard Schwarz. Online adaptation of language models with a memory of amortized contexts. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang, editors, *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024. http://papers.nips.cc/paper_files/paper/2024/hash/eaf956b52bae51fbf387b8be4cc3ce18-Abstract-Conference.html.

Haoran Tan, Zeyu Zhang, Chen Ma, Xu Chen, Quanyu Dai, and Zhenhua Dong. MemBench: Towards more comprehensive evaluation on the memory of LLM-based agents. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 19336–19352, Vienna, Austria, July 2025a. Association for Computational Linguistics. ISBN 979-8-89176-256-5. doi: 10.18653/v1/2025.findings-acl.989. https://aclanthology.org/2025.findings-acl.989/.

Xingyu Tan, Xiaoyang Wang, Xiwei Xu, Xin Yuan, Liming Zhu, and Wenjie Zhang. Memotime: Memory-augmented temporal knowledge graph enhanced large language model reasoning. *arXiv preprint arXiv:2510.13614*, 2025b.

Zhen Tan, Jun Yan, I-Hung Hsu, Rujun Han, Zifeng Wang, Long T. Le, Yiwen Song, Yanfei Chen, Hamid Palangi, George Lee, Anand Rajan Iyer, Tianlong Chen, Huan Liu, Chen-Yu Lee, and Tomas Pfister. In prospect and retrospect: Reflective memory management for long-term personalized dialogue agents. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar, editors, *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pages 8416–8439. Association for Computational Linguistics, 2025c. https://aclanthology.org/2025.acl-long.413/.

Hanlin Tang, Yang Lin, Jing Lin, Qingsen Han, Danning Ke, Shikuan Hong, Yiwu Yao, and Gongyi Wang. Razorattention: Efficient KV cache compression through retrieval heads. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net, 2025a. https://openreview.net/forum?id=tkiZQlL04w.

Qiaoyu Tang, Hao Xiang, Le Yu, Bowen Yu, Yaojie Lu, Xianpei Han, Le Sun, WenJuan Zhang, Pengbo Wang, Shixuan Liu, Zhenru Zhang, Jianhong Tu, Hongyu Lin, and Junyang Lin. Beyond turn limits: Training deep search agents with dynamic context window, October 2025b.

Xiangru Tang, Tianyu Hu, Muyang Ye, Yanjun Shao, Xunjian Yin, Siru Ouyang, Wangchunshu Zhou, Pan Lu, Zhuosheng Zhang, Yilun Zhao, Arman Cohan, and Mark Gerstein. Chemagent: Self-updating library in large language models improves chemical reasoning, 2025c. https://arxiv.org/abs/2501.06590.

Xiangru Tang, Tianrui Qin, Tianhao Peng, Ziyang Zhou, Daniel Shao, Tingting Du, Xinming Wei, Peng Xia, Fang Wu, He Zhu, Ge Zhang, Jiaheng Liu, Xingyao Wang, Sirui Hong, Chenglin Wu, Hao Cheng, Chi Wang, and Wangchunshu Zhou. Agent kb: Leveraging cross-domain experience for agentic problem solving, 2025d. https://arxiv.org/abs/2507.06229.

Yimin Tang, Yurong Xu, Ning Yan, and Masood Mortazavi. Enhancing long context performance in llms through inner loop query mechanism, 2024. https://arxiv.org/abs/2410.12859.

Dmitrii Tarasov, Elizaveta Goncharova, and Kuznetsov Andrey. Sentence-anchored gist compression for long-context llms, 2025. https://arxiv.org/abs/2511.08128.

Yi Tay, Vinh Tran, Mostafa Dehghani, Jianmo Ni, Dara Bahri, Harsh Mehta, Zhen Qin, Kai Hui, Zhe Zhao, Jai Prakash Gupta, Tal Schuster, William W. Cohen, and Donald Metzler. Transformer memory as a differentiable search index. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS*

*2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022. http://papers.nips.cc/paper_files/paper/2022/hash/892840a6123b5ec99ebaab8be1530fba-Abstract-Conference.html.

Shulin Tian, Ziniu Zhang, Liangyu Chen, and Ziwei Liu. Mmina: Benchmarking multihop multimodal internet agents, 2025. https://arxiv.org/abs/2404.09992.

Hieu Tran, Zonghai Yao, Nguyen Luong Tran, Zhichao Yang, Feiyun Ouyang, Shuo Han, Razieh Rahimi, and Hong Yu. PRIME: planning and retrieval-integrated memory for enhanced reasoning. *CoRR*, abs/2509.22315, 2025. doi: 10.48550/ARXIV.2509.22315. https://doi.org/10.48550/arXiv.2509.22315.

Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. Musique: Multihop questions via single-hop question composition, 2022. https://arxiv.org/abs/2108.00573.

Endel Tulving. Episodic and semantic memory. In *Organization of Memory*, pages xiii, 423–xiii, 423. Academic Press, Oxford, England, 1972.

Endel Tulving. Episodic memory: From mind to brain. *Annual Review of Psychology*, 53:1–25, 2002. ISSN 0066-4308. doi: 10.1146/annurev.psych.53.100901.135114.

Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. Zephyr: Direct distillation of lm alignment, 2023. https://arxiv.org/abs/2310.16944.

Szymon Tworkowski, Konrad Staniszewski, Mikolaj Pacek, Yuhuai Wu, Henryk Michalewski, and Piotr Milos. Focused transformer: Contrastive training for context scaling. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. http://papers.nips.cc/paper_files/paper/2023/hash/8511d06d5590f4bda24d42087802cc81-Abstract-Conference.html.

Luanbo Wan and Weizhi Ma. Storybench: A dynamic benchmark for evaluating long-term memory with multi turns. *arXiv preprint arXiv:2506.13356*, 2025.

Ziyu Wan, Yunxiang Li, Xiaoyu Wen, Yan Song, Hanjing Wang, Linyi Yang, Mark Schmidt, Jun Wang, Weinan Zhang, Shuyue Hu, and Ying Wen. Rema: Learning to meta-think for llms with multi-agent reinforcement learning, 2025. https://arxiv.org/abs/2503.09501.

Bing Wang, Xinnian Liang, Jian Yang, Hui Huang, Shuangzhi Wu, Peihao Wu, Lu Lu, Zejun Ma, and Zhoujun Li. Scm: Enhancing large language model with self-controlled memory framework, 2025a. https://arxiv.org/abs/2304.13343.

Bo Wang, Weiyi He, Shenglai Zeng, Zhen Xiang, Yue Xing, Jiliang Tang, and Pengfei He. Unveiling privacy risks in llm agent memory, 2025b. https://arxiv.org/abs/2502.13172.

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, 2025c.

Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *Trans. Mach. Learn. Res.*, 2024, 2024a. https://openreview.net/forum?id=ehfRiF0R3a.

Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An Open-Ended Embodied Agent with Large Language Models. *Trans. Mach. Learn. Res.*, 2024, 2024b. https://openreview.net/forum?id=ehfRiF0R3a.

Han Wang, Archiki Prasad, Elias Stengel-Eskin, and Mohit Bansal. Retrieval-augmented generation with conflicting evidence. *arXiv preprint arXiv:2504.13079*, 2025d.

Haochun Wang, Chi Liu, Nuwa Xi, Zewen Qiang, Sendong Zhao, Bing Qin, and Ting Liu. Huatuo: Tuning llama model with chinese medical knowledge, 2023a. https://arxiv.org/abs/2304.06975.

Hengyi Wang, Haizhou Shi, Shiwei Tan, Weiyi Qin, Wenyuan Wang, Tunyu Zhang, Akshay Nambi, Tanuja Ganu, and Hao Wang. Multimodal needle in a haystack: Benchmarking long-context capability of multimodal large language models. In Luis Chiruzzo, Alan Ritter, and Lu Wang, editors, *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 3221–3241, Albuquerque, New Mexico, April 2025e. Association for Computational Linguistics. ISBN 979-8-89176-189-6. doi: 10.18653/v1/2025.naacl-long.166. https://aclanthology.org/2025.naacl-long.166/.

Juyuan Wang, Rongchen Zhao, Wei Wei, Yufeng Wang, Mo Yu, Jie Zhou, Jin Xu, and Liyan Xu. Comorag: A cognitive-inspired memory-organized RAG for stateful long narrative reasoning. *CoRR*, abs/2508.10419, 2025f. doi: 10.48550/ARXIV.2508.10419. https://doi.org/10.48550/arXiv.2508.10419.

Mengru Wang, Ningyu Zhang, Ziwen Xu, Zekun Xi, Shumin Deng, Yunzhi Yao, Qishen Zhang, Linyi Yang, Jindong Wang, and Huajun Chen. Detoxifying large language models via knowledge editing. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 3093–3118. Association for Computational Linguistics, 2024c. doi: 10.18653/V1/2024.ACL-LONG.171. https://doi.org/10.18653/v1/2024.acl-long.171.

Noah Wang, Zhongyuan Peng, Haoran Que, Jiaheng Liu, Wangchunshu Zhou, Yuhan Wu, Hongcheng Guo, Ruitong Gan, Zehao Ni, Jian Yang, Man Zhang, Zhaoxiang Zhang, Wanli Ouyang, Ke Xu, Wenhao Huang, Jie Fu, and Junran Peng. Rolellm: Benchmarking, eliciting, and enhancing role-playing abilities of large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 14743–14777. Association for Computational Linguistics, 2024d. doi: 10.18653/V1/2024.FINDINGS-ACL.878. https://doi.org/10.18653/v1/2024.findings-acl.878.

Peng Wang, Zexi Li, Ningyu Zhang, Ziwen Xu, Yunzhi Yao, Yong Jiang, Pengjun Xie, Fei Huang, and Huajun Chen. WISE: rethinking the knowledge memory for lifelong model editing of large language models. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang, editors, *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024e. http://papers.nips.cc/paper_files/paper/2024/hash/60960ad78868fce5c165295fbd895060-Abstract-Conference.html.

Piaohong Wang, Motong Tian, Jiaxian Li, Yuan Liang, Yuqing Wang, Qianben Chen, Tiannan Wang, Zhicong Lu, Jiawei Ma, Yuchen Eleanor Jiang, and Wangchunshu Zhou. O-mem: Omni memory system for personalized, long horizon, self-evolving agents, 2025g. https://arxiv.org/abs/2511.13593.

Qingyue Wang, Yanhe Fu, Yanan Cao, Shuai Wang, Zhiliang Tian, and Liang Ding. Recursively summarizing enables long-term dialogue memory in large language models. *Neurocomputing*, 639:130193, 2025h.

Renxi Wang, Xudong Han, Lei Ji, Shu Wang, Timothy Baldwin, and Haonan Li. ToolGen: Unified tool retrieval and calling via generation. In *The Thirteenth International Conference on Learning Representations*, 2025i.

Ruize Wang, Duyu Tang, Nan Duan, Zhongyu Wei, Xuanjing Huang, Jianshu Ji, Guihong Cao, Daxin Jiang, and Ming Zhou. K-adapter: Infusing knowledge into pre-trained models with adapters. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Findings of the Association for Computational Linguistics: ACL/IJCNLP 2021, Online Event, August 1-6, 2021*, volume ACL/IJCNLP 2021 of *Findings of ACL*, pages 1405–1418. Association for Computational Linguistics, 2021. doi: 10.18653/V1/2021.FINDINGS-ACL.121. https://doi.org/10.18653/v1/2021.findings-acl.121.

Ruoyao Wang, Peter Jansen, Marc-Alexandre Côté, and Prithviraj Ammanabrolu. Scienceworld: Is your agent smarter than a 5th grader?, 2022a. https://arxiv.org/abs/2203.07540.

Tiannan Wang, Meiling Tao, Ruoyu Fang, Huilin Wang, Shuai Wang, Yuchen Eleanor Jiang, and Wangchunshu Zhou. Ai persona: Towards life-long personalization of llms, 2024f. https://arxiv.org/abs/2412.13103.

Weizhi Wang, Li Dong, Hao Cheng, Xiaodong Liu, Xifeng Yan, Jianfeng Gao, and Furu Wei. Augmenting language models with long-term memory. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023b. http://papers.nips.cc/paper_files/paper/2023/hash/ebd82705f44793b6f9ade5a669d0f0bf-Abstract-Conference.html.

Wenyi Wang, Piotr Piękos, Li Nanbo, Firas Laakom, Yimeng Chen, Mateusz Ostaszewski, Mingchen Zhuge, and Jürgen Schmidhuber. Huxley-godel machine: Human-level coding agent development by an approximation of the optimal self-improving machine, 2025j. https://arxiv.org/abs/2510.21614.

Xiaohan Wang, Yuhui Zhang, Orr Zohar, and Serena Yeung-Levy. VideoAgent: Long-Form Video Understanding with Large Language Model as Agent. In Ales Leonardis, Elisa Ricci, Stefan Roth, Olga Russakovsky, Torsten Sattler, and Gül Varol, editors, *Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part LXXX*, volume 15138 of *Lecture Notes in Computer Science*, pages 58–76. Springer, 2024g. doi: 10.1007/978-3-031-72989-8_4.

Yancheng Wang, Ziyan Jiang, Zheng Chen, Fan Yang, Yingxue Zhou, Eunah Cho, Xing Fan, Yanbin Lu, Xiaojiang Huang, and Yingzhen Yang. RecMind: Large Language Model Powered Agent For Recommendation. In Kevin Duh, Helena Gómez-Adorno, and Steven Bethard, editors, *Findings of the Association for Computational Linguistics: NAACL 2024, Mexico City, Mexico, June 16-21, 2024*, pages 4351–4364. Association for Computational Linguistics, 2024h. doi: 10.18653/V1/2024.FINDINGS-NAACL.271. https://doi.org/10.18653/v1/2024.findings-naacl.271.

Yanlin Wang, Wanjun Zhong, Yanxian Huang, Ensheng Shi, Min Yang, Jiachi Chen, Hui Li, Yuchi Ma, Qianxiang Wang, and Zibin Zheng. Agents in software engineering: Survey, landscape, and vision, 2024i. https://arxiv.org/abs/2409.09030.

Yingxu Wang, Siwei Liu, Jinyuan Fang, and Zaiqiao Meng. EvoAgentX: An automated framework for evolving agentic workflows. *arXiv preprint arXiv:2507.03616*, 2025k.

Yiping Wang, Qing Yang, Zhiyuan Zeng, Liliang Ren, Liyuan Liu, Baolin Peng, Hao Cheng, Xuehai He, Kuan Wang, Jianfeng Gao, Weizhu Chen, Shuohang Wang, Simon Shaolei Du, and Yelong Shen. Reinforcement learning for reasoning in large language models with one training example, 2025l. https://arxiv.org/abs/2504.20571.

Yu Wang and Xi Chen. MIRIX: Multi-agent memory system for llm-based agents. *arXiv preprint arXiv:2507.07957*, 2025.

Yu Wang, Xiusi Chen, Jingbo Shang, and Julian McAuley. Memoryllm: Towards self-updatable large language models. *ArXiv*, abs/2402.04624, 2024j. https://api.semanticscholar.org/CorpusID:267523037.

Yu Wang, Dmitry Krotov, Yuanzhe Hu, Yifan Gao, Wangchunshu Zhou, Julian J. McAuley, Dan Gutfreund, Rogério Feris, and Zexue He. M+: extending memoryllm with scalable long-term memory. *CoRR*, abs/2502.00592, 2025m. doi: 10.48550/ARXIV.2502.00592. https://doi.org/10.48550/arXiv.2502.00592.

Yu Wang, Xinshuang Liu, Xiusi Chen, Sean O'Brien, Junda Wu, and Julian J. McAuley. Self-updatable large language models by integrating context into model parameters. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net, 2025n. https://openreview.net/forum?id=aCPFCDL9QY.

Yu Wang, Ryuichi Takanobu, Zhiqi Liang, Yuzhen Mao, Yuanzhe Hu, Julian J. McAuley, and Xiaojian Wu. Mem-$\alpha$: Learning memory construction via reinforcement learning. *CoRR*, abs/2509.25911, 2025o. doi: 10.48550/ARXIV. 2509.25911. https://doi.org/10.48550/arXiv.2509.25911.

Yujing Wang, Yingyan Hou, Haonan Wang, Ziming Miao, Shibin Wu, Qi Chen, Yuqing Xia, Chengmin Chi, Guoshuai Zhao, Zheng Liu, Xing Xie, Hao Sun, Weiwei Deng, Qi Zhang, and Mao Yang. A neural corpus indexer for document retrieval. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022b. http://papers.nips.cc/paper_files/paper/2022/hash/a46156bd3579c3b268108ea6aca71d13-Abstract-Conference.html.

Zheng Wang, Zhongyang Li, Zeren Jiang, Dandan Tu, and Wei Shi. Crafting personalized agents through retrieval-augmented generation on editable memory graphs, 2024k. https://arxiv.org/abs/2409.19401.

Zihao Wang, Shaofei Cai, Anji Liu, Yonggang Jin, Jinbing Hou, Bowei Zhang, Haowei Lin, Zhaofeng He, Zilong Zheng, Yaodong Yang, Xiaojian Ma, and Yitao Liang. JARVIS-1: open-world multi-task agents with memory-augmented multimodal language models. *IEEE Trans. Pattern Anal. Mach. Intell.*, 47(3):1894–1907, 2025p. doi: 10.1109/TPAMI.2024.3511593. https://doi.org/10.1109/TPAMI.2024.3511593.

Zixuan Wang, Bo Yu, Junzhe Zhao, Wenhao Sun, Sai Hou, Shuai Liang, Xing Hu, Yinhe Han, and Yiming Gan. KARMA: Augmenting Embodied AI Agents with Long-and-Short Term Memory Systems. In *IEEE International Conference on Robotics and Automation, ICRA 2025, Atlanta, GA, USA, May 19-23, 2025*, pages 1–8. IEEE, 2025q. doi: 10.1109/ICRA55743.2025.11128047.

Zora Zhiruo Wang, Jiayuan Mao, Daniel Fried, and Graham Neubig. Agent workflow memory. In *Forty-second International Conference on Machine Learning*, 2024l.

Zora Zhiruo Wang, Apurva Gandhi, Graham Neubig, and Daniel Fried. Inducing programmatic skills for agentic tasks, 2025r. https://arxiv.org/abs/2504.06821.

Joel Ward. Memoriesdb: A temporal-semantic-relational database for long-term agent memory/modeling experience as a graph of temporal-semantic surfaces. *arXiv preprint arXiv:2511.06179*, 2025.

Christopher Watkins and Peter Dayan. Q-learning. *Machine Learning*, 8:279–292, 1992. https://api.semanticscholar.org/CorpusID:208910339.

Haoran Wei, Yaofeng Sun, and Yukun Li. Deepseek-ocr: Contexts optical compression. *arXiv preprint arXiv:2510.18234*, 2025a.

Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. Finetuned language models are zero-shot learners, 2022. https://arxiv.org/abs/2109.01652.

Jason Wei, Zhiqing Sun, Spencer Papay, Scott McKinney, Jeffrey Han, Isa Fulford, Hyung Won Chung, Alex Tachard Passos, William Fedus, and Amelia Glaese. BrowseComp: A simple yet challenging benchmark for browsing agents. *arXiv preprint arXiv:2504.12516*, 2025b.

Jiaqi Wei, Yuejin Yang, Xiang Zhang, Yuhan Chen, Xiang Zhuang, Zhangyang Gao, Dongzhan Zhou, Guangshuai Wang, Zhiqiang Gao, Juntai Cao, Zijie Qiu, Ming Hu, Chenglong Ma, Shixiang Tang, Junjun He, Chunfeng Song, Xuming He, Qiang Zhang, Chenyu You, Shuangjia Zheng, Ning Ding, Wanli Ouyang, Nanqing Dong, Yu Cheng, Siqi Sun, Lei Bai, and Bowen Zhou. From ai for science to agentic science: A survey on autonomous scientific discovery, 2025c. https://arxiv.org/abs/2508.14111.

Rubin Wei, Jiaqi Cao, Jiarui Wang, Jushi Kai, Qipeng Guo, Bowen Zhou, and Zhouhan Lin. MLP memory: Language modeling with retriever-pretrained external memory. *CoRR*, abs/2508.01832, 2025d. doi: 10.48550/ARXIV.2508.01832. https://doi.org/10.48550/arXiv.2508.01832.

Tianxin Wei, Noveen Sachdeva, Benjamin Coleman, Zhankui He, Yuanchen Bei, Xuying Ning, Mengting Ai, Yunzhe Li, Jingrui He, Ed H Chi, et al. Evo-memory: Benchmarking llm agent test-time learning with self-evolving memory. *arXiv preprint arXiv:2511.20857*, 2025e.

Yixuan Weng, Minjun Zhu, Qiujie Xie, Qiyao Sun, Zhen Lin, Sifan Liu, and Yue Zhang. Deepscientist: Advancing frontier-pushing scientific findings progressively, 2025. https://arxiv.org/abs/2509.26603.

Rebecca Westhäußer, Wolfgang Minker, and Sebastian Zepf. Enabling personalized long-term interactions in llm-based agents through persistent memory and user profiles. *CoRR*, abs/2510.07925, 2025. doi: 10.48550/ARXIV.2510.07925. https://doi.org/10.48550/arXiv.2510.07925.

Martin Wistuba, Prabhu Teja Sivaprasad, Lukas Balles, and Giovanni Zappella. Continual learning with low rank adaptation. *CoRR*, abs/2311.17601, 2023. doi: 10.48550/ARXIV.2311.17601. https://doi.org/10.48550/arXiv.2311.17601.

Cheng-Kuang Wu, Zhi Rui Tam, Chieh-Yen Lin, Yun-Nung Chen, and Hung-yi Lee. Streambench: Towards benchmarking continuous improvement of language agents. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang, editors, *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024a. http://papers.nips.cc/paper_files/paper/2024/hash/c189915371c4474fe9789be3728113fc-Abstract-Datasets_and_Benchmarks_Track.html.

Di Wu, Hongwei Wang, Wenhao Yu, Yuwei Zhang, Kai-Wei Chang, and Dong Yu. Longmemeval: Benchmarking chat assistants on long-term interactive memory. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net, 2025a. https://openreview.net/forum?id=pZiyCaVuti.

Jeff Wu, Long Ouyang, Daniel M Ziegler, Nisan Stiennon, Ryan Lowe, Jan Leike, and Paul Christiano. Recursively summarizing books with human feedback. *arXiv preprint arXiv:2109.10862*, 2021.

Jialong Wu, Baixuan Li, Runnan Fang, Wenbiao Yin, Liwen Zhang, Zhengwei Tao, Dingchu Zhang, Zekun Xi, Yong Jiang, Pengjun Xie, Fei Huang, and Jingren Zhou. Webdancer: Towards autonomous information seeking agency. *CoRR*, abs/2505.22648, 2025b. doi: 10.48550/ARXIV.2505.22648. https://doi.org/10.48550/arXiv.2505.22648.

Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, et al. Autogen: Enabling next-gen llm applications via multi-agent conversations. In *First Conference on Language Modeling*, 2024b.

Rong Wu, Xiaoman Wang, Jianbiao Mei, Pinlong Cai, Daocheng Fu, Cheng Yang, Licheng Wen, Xuemeng Yang, Yufan Shen, Yuxin Wang, and Botian Shi. Evolver: Self-evolving llm agents through an experience-driven lifecycle, 2025c. https://arxiv.org/abs/2510.16079.

Wenyi Wu, Zixuan Song, Kun Zhou, Yifei Shao, Zhiting Hu, and Biwei Huang. Towards general continuous memory for vision-language models. *ArXiv*, abs/2505.17670, 2025d.

Wenyi Wu, Kun Zhou, Ruoxin Yuan, Vivian Yu, Stephen Wang, Zhiting Hu, and Biwei Huang. Auto-scaling continuous memory for GUI agent. *CoRR*, abs/2510.09038, 2025e. doi: 10.48550/ARXIV.2510.09038. https://doi.org/10.48550/arXiv.2510.09038.

Xixi Wu, Kuan Li, Yida Zhao, Liwen Zhang, Litu Ou, Huifeng Yin, Zhongwang Zhang, Yong Jiang, Pengjun Xie, Fei Huang, Minhao Cheng, Shuai Wang, Hong Cheng, and Jingren Zhou. ReSum: Unlocking Long-Horizon Search Intelligence via Context Summarization. *CoRR*, abs/2509.13313, 2025f. doi: 10.48550/ARXIV.2509.13313.

Yaxiong Wu, Sheng Liang, Chen Zhang, Yichao Wang, Yongyue Zhang, Huifeng Guo, Ruiming Tang, and Yong Liu. From human memory to ai memory: A survey on memory mechanisms in the era of llms, 2025g. https://arxiv.org/abs/2504.15965.

Yaxiong Wu, Yongyue Zhang, Sheng Liang, and Yong Liu. Sgmem: Sentence graph memory for long-term conversational agents. *ArXiv*, abs/2509.21212, 2025h.

Yiran Wu, Feiran Jia, Shaokun Zhang, Hangyu Li, Erkang Zhu, Yue Wang, Yin Tat Lee, Richard Peng, Qingyun Wu, and Chi Wang. Mathchat: Converse to tackle challenging math problems with llm agents, 2024c. https://arxiv.org/abs/2306.01337.

Yisha Wu, Cen Zhao, Yuanpei Cao, Xiaoqing Xu, Yashar Mehdad, Mindy Ji, and Claire Na Cheng. Incremental summarization for customer support via progressive note-taking and agent feedback. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 2000–2015, 2025i.

Yuhuai Wu, Markus Norman Rabe, DeLesley Hutchins, and Christian Szegedy. Memorizing transformers. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. https://openreview.net/forum?id=TrjbxzRcnf-.

Zijun Wu, Yongchang Hao, and Lili Mou. Tokmem: Tokenized procedural memory for large language models, 2025j. https://arxiv.org/abs/2510.00444.

Rui Xi and Xianghan Wang. Livia: An emotion-aware AR companion powered by modular AI agents and progressive memory compression. *CoRR*, abs/2509.05298, 2025. doi: 10.48550/ARXIV.2509.05298. https://doi.org/10.48550/arXiv.2509.05298.

Yunjia Xi, Weiwen Liu, Jianghao Lin, Bo Chen, Ruiming Tang, Weinan Zhang, and Yong Yu. Memocrs: Memory-enhanced sequential conversational recommender systems with large language models. In Edoardo Serra and Francesca Spezzano, editors, *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management, CIKM 2024, Boise, ID, USA, October 21-25, 2024*, pages 2585–2595. ACM, 2024a. doi: 10.1145/3627673.3679599. https://doi.org/10.1145/3627673.3679599.

Zhiheng Xi, Yiwen Ding, Wenxiang Chen, Boyang Hong, Honglin Guo, Junzhe Wang, Dingwen Yang, Chenyang Liao, Xin Guo, Wei He, Songyang Gao, Lu Chen, Rui Zheng, Yicheng Zou, Tao Gui, Qi Zhang, Xipeng Qiu, Xuanjing Huang, Zuxuan Wu, and Yu-Gang Jiang. Agentgym: Evolving large language model-based agents across diverse environments, 2024b. https://arxiv.org/abs/2406.04151.

Siyu Xia, Zekun Xu, Jiajun Chai, Wentian Fan, Yan Song, Xiaohan Wang, Guojun Yin, Wei Lin, Haifeng Zhang, and Jun Wang. From experience to strategy: Empowering llm agents with trainable graph memory. *arXiv preprint arXiv:2511.07800*, 2025.

Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming language models with attention sinks. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. https://openreview.net/forum?id=NG7sS51zVF.

Yuan-An Xiao, Pengfei Gao, Chao Peng, and Yingfei Xiong. Improving the efficiency of llm agent systems through trajectory reduction, 2025a. https://arxiv.org/abs/2509.23586.

Yunzhong Xiao, Yangmin Li, Hewei Wang, Yunlong Tang, and Zora Zhiruo Wang. ToolMem: Enhancing Multimodal Agents with Learnable Tool Capability Memory. *CoRR*, abs/2510.06664, 2025b. doi: 10.48550/ARXIV.2510.06664. https://doi.org/10.48550/arXiv.2510.06664. arXiv: 2510.06664.

Zeqi Xiao, Yushi Lan, Yifan Zhou, Wenqi Ouyang, Shuai Yang, Yanhong Zeng, and Xingang Pan. Worldmem: Long-term consistent world simulation with memory. *arXiv preprint arXiv:2504.12369*, 2025c.

Junlin Xie, Zhihong Chen, Ruifei Zhang, Xiang Wan, and Guanbin Li. Large multimodal agents: A survey, 2024. https://arxiv.org/abs/2402.15116.

Haoran Xu, Jiacong Hu, Ke Zhang, Lei Yu, Yuxin Tang, Xinyuan Song, Yiqun Duan, Lynn Ai, and Bill Shi. Sedm: Scalable self-evolving distributed memory for agents, 2025a. https://arxiv.org/abs/2509.09498.

Mufan Xu, Gewen Liang, Kehai Chen, Wei Wang, Xun Zhou, Muyun Yang, Tiejun Zhao, and Min Zhang. Memory-augmented query reconstruction for llm-based knowledge graph reasoning, 2025b. https://arxiv.org/abs/2503.05193.

Renjun Xu and Jingwen Peng. A comprehensive survey of deep research: Systems, methodologies, and applications, 2025. https://arxiv.org/abs/2506.12594.

Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao, Juntao Tan, and Yongfeng Zhang. A-MEM: Agentic Memory for LLM Agents. *CoRR*, abs/2502.12110, 2025c. doi: 10.48550/ARXIV.2502.12110.

Yige Xu, Xu Guo, Zhiwei Zeng, and Chunyan Miao. Softcot: Soft chain-of-thought for efficient reasoning with llms. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar, editors, *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pages 23336–23351. Association for Computational Linguistics, 2025d. https://aclanthology.org/2025.acl-long.1137/.

B. Y. Yan, Chaofan Li, Hongjin Qian, Shuqi Lu, and Zheng Liu. General agentic memory via deep research, November 2025a.

Ming Yan, Ruihao Li, Hao Zhang, Hao Wang, Zhilan Yang, and Ji Yan. Larp: Language-agent role play for open-world games, 2023. https://arxiv.org/abs/2312.17653.

Sikuan Yan, Xiufeng Yang, Zuchao Huang, Ercong Nie, Zifeng Ding, Zonggen Li, Xiaowen Ma, Hinrich Schütze, Volker Tresp, and Yunpu Ma. Memory-R1: Enhancing large language model agents to manage and utilize memories via reinforcement learning. *arXiv preprint arXiv:2508.19828*, 2025b.

Hongkang Yang, Zehao Lin, Wenjin Wang, Hao Wu, Zhiyu Li, Bo Tang, Wenqiang Wei, Jinbo Wang, Zeyun Tang, Shichao Song, Chenyang Xi, Yu Yu, Kai Chen, Feiyu Xiong, Linpeng Tang, and Weinan E. Memory[3]: Language modeling with explicit memory. *CoRR*, abs/2407.01178, 2024a. doi: 10.48550/ARXIV.2407.01178. https://doi.org/10.48550/arXiv.2407.01178.

Ling Yang, Zhaochen Yu, Tianjun Zhang, Shiyi Cao, Minkai Xu, Wentao Zhang, Joseph E. Gonzalez, and Bin Cui. Buffer of thoughts: Thought-augmented reasoning with large language models. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang, editors, *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024b. http://papers.nips.cc/paper_files/paper/2024/hash/cde328b7bf6358f5ebb91fe9c539745e-Abstract-Conference.html.

Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering, 2018. https://arxiv.org/abs/1809.09600.

Ziyi Yang, Zaibin Zhang, Zirui Zheng, Yuxian Jiang, Ziyue Gan, Zhiyu Wang, Zijian Ling, Jinsong Chen, Martz Ma, Bowen Dong, Prateek Gupta, Shuyue Hu, Zhenfei Yin, Guohao Li, Xu Jia, Lijun Wang, Bernard Ghanem, Huchuan Lu, Chaochao Lu, Wanli Ouyang, Yu Qiao, Philip Torr, and Jing Shao. Oasis: Open agent social interaction simulations with one million agents, 2025. https://arxiv.org/abs/2411.11581.

Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents, 2023a. https://arxiv.org/abs/2207.01206.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan Cao. ReAct: Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations*, 2023b.

Weiran Yao, Shelby Heinecke, Juan Carlos Niebles, Zhiwei Liu, Yihao Feng, Le Xue, Rithesh R. N., Zeyuan Chen, Jianguo Zhang, Devansh Arpit, Ran Xu, Phil Mui, Huan Wang, Caiming Xiong, and Silvio Savarese. Retroformer: Retrospective large language agents with policy gradient optimization. In *The Twelfth International Conference on Learning Representations*, 2024a.

Yao Yao, Zuchao Li, and Hai Zhao. Sirllm: Streaming infinite retentive LLM. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 2611–2624. Association for Computational Linguistics, 2024b. doi: 10.18653/V1/2024.ACL-LONG.143. https://doi.org/10.18653/v1/2024.acl-long.143.

Rui Ye, Zhongwang Zhang, Kuan Li, Huifeng Yin, Zhengwei Tao, Yida Zhao, Liangcai Su, Liwen Zhang, Zile Qiao, Xinyu Wang, et al. Agentfold: Long-horizon web agents with proactive context management. *arXiv preprint arXiv:2510.24699*, 2025a.

Shicheng Ye, Chao Yu, Kaiqiang Ke, Chengdong Xu, and Yinqi Wei. $H^2$r: Hierarchical hindsight reflection for multi-task LLM agents. *CoRR*, abs/2509.12810, 2025b. doi: 10.48550/ARXIV.2509.12810. https://doi.org/10.48550/arXiv.2509.12810.

Ryan Yen and Jian Zhao. Memolet: Reifying the Reuse of User-AI Conversational Memories. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology*, UIST '24, pages 1–22, New York, NY, USA, October 2024. Association for Computing Machinery. ISBN 979-8-4007-0628-8. doi: 10.1145/3654777.3676388.

Xunjian Yin, Xinyi Wang, Liangming Pan, Li Lin, Xiaojun Wan, and William Yang Wang. Gödel agent: A self-referential agent framework for recursively self-improvement. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar, editors, *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pages 27890–27913. Association for Computational Linguistics, 2025. https://aclanthology.org/2025.acl-long.1354/.

Chanwoong Yoon, Taewhoo Lee, Hyeon Hwang, Minbyul Jeong, and Jaewoo Kang. CompAct: Compressing Retrieved Documents Actively for Question Answering. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 21424–21439. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.EMNLP-MAIN.1194.

Zeng You, Zhiquan Wen, Yaofo Chen, Xin Li, Runhao Zeng, Yaowei Wang, and Mingkui Tan. Towards long video understanding via fine-detailed video story generation. *IEEE Transactions on Circuits and Systems for Video Technology*, 2024.

Hongli Yu, Tinghong Chen, Jiangtao Feng, Jiangjie Chen, Weinan Dai, Qiying Yu, Ya-Qin Zhang, Wei-Ying Ma, Jingjing Liu, Mingxuan Wang, et al. Memagent: Reshaping long-context llm with multi-conv rl-based memory agent. *arXiv preprint arXiv:2507.02259*, 2025a.

Jiwen Yu, Jianhong Bai, Yiran Qin, Quande Liu, Xintao Wang, Pengfei Wan, Di Zhang, and Xihui Liu. Context as memory: Scene-consistent interactive long video generation with memory retrieval. *CoRR*, abs/2506.03141, 2025b. doi: 10.48550/ARXIV.2506.03141. https://doi.org/10.48550/arXiv.2506.03141.

Jiwen Yu, Jianhong Bai, Yiran Qin, Quande Liu, Xintao Wang, Pengfei Wan, Di Zhang, and Xihui Liu. Context as memory: Scene-consistent interactive long video generation with memory retrieval. *arXiv preprint arXiv:2506.03141*, 2025c.

Tao Yu, Zhengbo Zhang, Zhiheng Lyu, Junhao Gong, Hongzhu Yi, Xinming Wang, Yuxuan Zhou, Jiabing Yang, Ping Nie, Yan Huang, and Wenhu Chen. BrowserAgent: Building Web Agents with Human-Inspired Web Browsing Actions, October 2025d. http://arxiv.org/abs/2510.10666. arXiv:2510.10666 [cs].

Xinlei Yu, Chengming Xu, Guibin Zhang, Zhangquan Chen, Yudong Zhang, Yongbo He, Peng-Tao Jiang, Jiangning Zhang, Xiaobin Hu, and Shuicheng Yan. Vismem: Latent vision memory unlocks potential of vision-language models, 2025e. https://arxiv.org/abs/2511.11007.

Yangyang Yu, Zhiyuan Yao, Haohang Li, Zhiyang Deng, Yuechen Jiang, Yupeng Cao, Zhi Chen, Jordan Suchow, Zhenyu Cui, Rong Liu, et al. Fincon: A synthesized llm multi-agent system with conceptual verbal reinforcement for enhanced financial decision making. *Advances in Neural Information Processing Systems*, 37:137010–137045, 2024.

Yifei Yu, Xiaoshan Wu, Xinting Hu, Tao Hu, Yangtian Sun, Xiaoyang Lyu, Bo Wang, Lin Ma, Yuewen Ma, Zhongrui Wang, et al. Videossm: Autoregressive long video generation with hybrid state-space memory. *arXiv preprint arXiv:2512.04519*, 2025f.

Lifan Yuan, Yangyi Chen, Xingyao Wang, Yi Fung, Hao Peng, and Heng Ji. CRAFT: Customizing LLMs by Creating and Retrieving from Specialized Toolsets. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024a. https://openreview.net/forum?id=G0vdDSt9XM.

Peiwen Yuan, Xinglin Wang, Shaoxiong Feng, Boyuan Pan, Yiwei Li, Heda Wang, Xupeng Miao, and Kan Li. Generative dense retrieval: Memory can be a burden. In Yvette Graham and Matthew Purver, editors, *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2024 - Volume 1: Long Papers, St. Julian's, Malta, March 17-22, 2024*, pages 2835–2845. Association for Computational Linguistics, 2024b. https://aclanthology.org/2024.eacl-long.173.

Qianhao Yuan, Jie Lou, Zichao Li, Jiawei Chen, Yaojie Lu, Hongyu Lin, Le Sun, Debing Zhang, and Xianpei Han. MemSearcher: Training LLMs to Reason, Search and Manage Memory via End-to-End Reinforcement Learning, November 2025a.

Ruifeng Yuan, Shichao Sun, Yongqi Li, Zili Wang, Ziqiang Cao, and Wenjie Li. Personalized large language model assistant with evolving conditional memory. In Owen Rambow, Leo Wanner, Marianna Apidianaki, Hend Al-Khalifa, Barbara Di Eugenio, and Steven Schockaert, editors, *Proceedings of the 31st International Conference on Computational Linguistics, COLING 2025, Abu Dhabi, UAE, January 19-24, 2025*, pages 3764–3777. Association for Computational Linguistics, 2025b. https://aclanthology.org/2025.coling-main.254/.

Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Xian Li, Sainbayar Sukhbaatar, Jing Xu, and Jason Weston. Self-rewarding language models. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024c. https://openreview.net/forum?id=0NphYCmgua.

Sizhe Yuen, Francisco Gomez Medina, Ting Su, Yali Du, and Adam J. Sobey. Intrinsic memory agents: Heterogeneous multi-agent llm systems through structured contextual memory, 2025. https://arxiv.org/abs/2508.08997.

Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontañón, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. Big bird: Transformers for longer sequences. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. https://proceedings.neurips.cc/paper/2020/hash/c8512d142a2d849725f31a9a7a361ab9-Abstract.html.

Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah Goodman. STaR: Bootstrapping reasoning with reasoning. In *Advances in Neural Information Processing Systems*, volume 35, pages 15476–15488, 2022.

Hansi Zeng, Chen Luo, Bowen Jin, Sheikh Muhammad Sarwar, Tianxin Wei, and Hamed Zamani. Scalable and effective generative information retrieval. In Tat-Seng Chua, Chong-Wah Ngo, Ravi Kumar, Hady W. Lauw, and Roy Ka-Wei Lee, editors, *Proceedings of the ACM on Web Conference 2024, WWW 2024, Singapore, May 13-17, 2024*, pages 1441–1452. ACM, 2024. doi: 10.1145/3589334.3645477. https://doi.org/10.1145/3589334.3645477.

Yunpeng Zhai, Shuchang Tao, Cheng Chen, Anni Zou, Ziqian Chen, Qingxu Fu, Shinji Mai, Li Yu, Jiaji Deng, Zouying Cao, Zhaoyang Liu, Bolin Ding, and Jingren Zhou. Agentevolver: Towards efficient self-evolving agent system, 2025. https://arxiv.org/abs/2511.10395.

Chaoyun Zhang, He Huang, Chiming Ni, Jian Mu, Si Qin, Shilin He, Lu Wang, Fangkai Yang, Pu Zhao, Chao Du, Liqun Li, Yu Kang, Zhao Jiang, Suzhen Zheng, Rujia Wang, Jiaxu Qian, Minghua Ma, Jian-Guang Lou, Qingwei Lin, Saravan Rajmohan, and Dongmei Zhang. UFO2: the desktop agentos. *CoRR*, abs/2504.14603, 2025a. doi: 10.48550/ARXIV.2504.14603. https://doi.org/10.48550/arXiv.2504.14603.

Gaoke Zhang, Bo Wang, Yunlong Ma, Dongming Zhao, and Zifei Yu. Multiple memory systems for enhancing the long-term memory of agent. *CoRR*, abs/2508.15294, 2025b. doi: 10.48550/ARXIV.2508.15294. https://doi.org/10.48550/arXiv.2508.15294.

Gui-Min Zhang, Muxin Fu, Guancheng Wan, Miao Yu, Kun Wang, and Shuicheng Yan. G-memory: Tracing hierarchical memory for multi-agent systems. *ArXiv*, abs/2506.07398, 2025c.

Gui-Min Zhang, Muxin Fu, and Shuicheng Yan. Memgen: Weaving generative latent memory for self-evolving agents. *ArXiv*, abs/2509.24704, 2025d.

Guibin Zhang, Muxin Fu, Guancheng Wan, Miao Yu, Kun Wang, and Shuicheng Yan. G-Memory: Tracing hierarchical memory for multi-agent systems. *arXiv preprint arXiv:2506.07398*, 2025e.

Guibin Zhang, Hejia Geng, Xiaohang Yu, Zhenfei Yin, Zaibin Zhang, Zelin Tan, Heng Zhou, Zhongzhi Li, Xiangyuan Xue, Yijiang Li, Yifan Zhou, Yang Chen, Chen Zhang, Yutao Fan, Zihu Wang, Songtao Huang, Yue Liao, Hongru Wang, Mengyue Yang, Heng Ji, Michael Littman, Jun Wang, Shuicheng Yan, Philip Torr, and Lei Bai. The landscape of agentic reinforcement learning for llms: A survey, 2025f. https://arxiv.org/abs/2509.02547.

Guibin Zhang, Fanci Meng, Guancheng Wan, Zherui Li, Kun Wang, Zhenfei Yin, Lei Bai, and Shuicheng Yan. Latentevolve: Self-evolving test-time scaling in latent space, 2025g. https://arxiv.org/abs/2509.24771.

Jenny Zhang, Shengran Hu, Cong Lu, Robert T. Lange, and Jeff Clune. Darwin Godel Machine: Open-Ended Evolution of Self-Improving Agents. *CoRR*, abs/2505.22954, 2025h. doi: 10.48550/ARXIV.2505.22954. https://doi.org/10.48550/arXiv.2505.22954. arXiv: 2505.22954.

Jiarui Zhang. Guided profile generation improves personalization with llms. *arXiv preprint arXiv:2409.13093*, 2024.

Jiayi Zhang, Jinyu Xiang, Zhaoyang Yu, Fengwei Teng, Xiong-Hui Chen, Jiaqi Chen, Mingchen Zhuge, Xin Cheng, Sirui Hong, Jinlin Wang, Bingnan Zheng, Bang Liu, Yuyu Luo, and Chenglin Wu. AFlow: Automating agentic workflow generation. In *The Thirteenth International Conference on Learning Representations*, 2025i.

Kai Zhang, Xiangchao Chen, Bo Liu, Tianci Xue, Zeyi Liao, Zhihan Liu, Xiyao Wang, Yuting Ning, Zhaorun Chen, Xiaohan Fu, Jian Xie, Yuxuan Sun, Boyu Gou, Qi Qi, Zihang Meng, Jianwei Yang, Ning Zhang, Xian Li, Ashish Shah, Dat Huynh, Hengduo Li, Zi Yang, Sara Cao, Lawrence Jang, Shuyan Zhou, Jiacheng Zhu, Huan Sun, Jason Weston, Yu Su, and Yifan Wu. Agent learning via early experience, 2025j. https://arxiv.org/abs/2510.08558.

Kaiyan Zhang, Yuxin Zuo, Bingxiang He, Youbang Sun, Runze Liu, Che Jiang, Yuchen Fan, Kai Tian, Guoli Jia, Pengfei Li, Yu Fu, Xingtai Lv, Yuchen Zhang, Sihang Zeng, Shang Qu, Haozhan Li, Shijie Wang, Yuru Wang, Xinwei Long, Fangfu Liu, Xiang Xu, Jiaze Ma, Xuekai Zhu, Ermo Hua, Yihao Liu, Zonglin Li, Huayu Chen, Xiaoye Qu, Yafu Li, Weize Chen, Zhenzhao Yuan, Junqi Gao, Dong Li, Zhiyuan Ma, Ganqu Cui, Zhiyuan Liu, Biqing Qi, Ning Ding, and Bowen Zhou. A survey of reinforcement learning for large reasoning models, 2025k. https://arxiv.org/abs/2509.08827.

Lingfeng Zhang, Yuecheng Liu, Zhanguang Zhang, Matin Aghaei, Yaochen Hu, Hongjian Gu, Mohammad Ali Alomrani, David Gamaliel Arcos Bravo, Raika Karimi, Atia Hamidizadeh, Haoping Xu, Guowei Huang, Zhanpeng Zhang, Tongtong Cao, Weichao Qiu, Xingyue Quan, Jianye Hao, Yuzheng Zhuang, and Yingxue Zhang. Mem2ego: Empowering vision-language models with global-to-ego memory for long-horizon embodied navigation. *CoRR*, abs/2502.14254, 2025l. doi: 10.48550/ARXIV.2502.14254. https://doi.org/10.48550/arXiv.2502.14254.

Qizheng Zhang, Changran Hu, Shubhangi Upasani, Boyuan Ma, Fenglu Hong, Vamsidhar Kamanuru, Jay Rainton, Chen Wu, Mengmeng Ji, Hanchen Li, Urmish Thakker, James Zou, and Kunle Olukotun. Agentic context engineering: Evolving contexts for self-improving language models. *CoRR*, abs/2510.04618, 2025m. doi: 10.48550/ARXIV.2510.04618. https://doi.org/10.48550/arXiv.2510.04618.

Shaohua Zhang, Yuan Lin, and Hang Li. Memory retrieval and consolidation in large language models through function tokens, 2025n. https://arxiv.org/abs/2510.08203.

Wenlin Zhang, Xiaopeng Li, Yingyi Zhang, Pengyue Jia, Yichao Wang, Huifeng Guo, Yong Liu, and Xiangyu Zhao. Deep research: A survey of autonomous research agents, 2025o. https://arxiv.org/abs/2508.12752.

Wentao Zhang, Lingxuan Zhao, Haochong Xia, Shuo Sun, Jiaze Sun, Molei Qin, Xinyi Li, Yuqing Zhao, Yilei Zhao, Xinyu Cai, Longtao Zheng, Xinrun Wang, and Bo An. A multimodal foundation agent for financial trading: Tool-augmented, diversified, and generalist, 2024. https://arxiv.org/abs/2402.18485.

Yaoze Zhang, Rong Wu, Pinlong Cai, Xiaoman Wang, Guohang Yan, Song Mao, Ding Wang, and Botian Shi. Leanrag: Knowledge-graph-based generation with semantic aggregation and hierarchical retrieval, 2025p. https://arxiv.org/abs/2508.10391.

Yuxiang Zhang, Jiangming Shu, Ye Ma, Xueyuan Lin, Shangxi Wu, and Jitao Sang. Memory as action: Autonomous context curation for long-horizon agentic tasks. *CoRR*, abs/2510.12635, 2025q. doi: 10.48550/ARXIV.2510.12635. https://doi.org/10.48550/arXiv.2510.12635.

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, 2025r.

Zeyu Zhang, Quanyu Dai, Xu Chen, Rui Li, Zhongyang Li, and Zhenhua Dong. Memengine: A unified and modular library for developing advanced memory of llm-based agents, 2025s. https://arxiv.org/abs/2505.02099.

Zeyu Zhang, Quanyu Dai, Rui Li, Xiaohe Bo, Xu Chen, and Zhenhua Dong. Learn to memorize: Optimizing llm-based agents with adaptive memory framework. *CoRR*, abs/2508.16629, 2025t. doi: 10.48550/ARXIV.2508.16629. https://doi.org/10.48550/arXiv.2508.16629.

Zeyu Zhang, Yang Zhang, Haoran Tan, Rui Li, and Xu Chen. Explicit vs implicit memory: Exploring multi-hop complex reasoning over personalized information. *arXiv preprint arXiv:2508.13250*, 2025u.

Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song, Yuandong Tian, Christopher Ré, Clark W. Barrett, Zhangyang Wang, and Beidi Chen. H2O: heavy-hitter oracle for efficient generative inference of large language models. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. http://papers.nips.cc/paper_files/paper/2023/hash/6ceefa7b15572587b78ecfcebb2827f8-Abstract-Conference.html.

Andrew Zhao, Daniel Huang, Quentin Xu, Matthieu Lin, Yong-Jin Liu, and Gao Huang. Expel: LLM agents are experiential learners. In Michael J. Wooldridge, Jennifer G. Dy, and Sriraam Natarajan, editors, *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014,*

*February 20-27, 2024, Vancouver, Canada*, pages 19632–19642. AAAI Press, 2024. doi: 10.1609/AAAI.V38I17.29936. https://doi.org/10.1609/aaai.v38i17.29936.

Di Zhao, Longhui Ma, Siwei Wang, Miao Wang, and Zhao Lv. COLA: A Scalable Multi-Agent Framework For Windows UI Task Automation. *CoRR*, abs/2503.09263, 2025a. doi: 10.48550/ARXIV.2503.09263. https://doi.org/10.48550/arXiv.2503.09263. arXiv: 2503.09263.

Linxi Zhao, Sofian Zalouk, Christian K. Belardi, Justin Lovelace, Jin Peng Zhou, Ryan Thomas Noonan, Dongyoung Go, Kilian Q. Weinberger, Yoav Artzi, and Jennifer J. Sun. Pre-training limited memory language models with internal and external knowledge, 2025b. https://arxiv.org/abs/2505.15962.

Shitian Zhao, Haoquan Zhang, Shaoheng Lin, Ming Li, Qilong Wu, Kaipeng Zhang, and Chen Wei. Pyvision: Agentic vision with dynamic tooling, 2025c. https://arxiv.org/abs/2507.07998.

Siyan Zhao, Mingyi Hong, Yang Liu, Devamanyu Hazarika, and Kaixiang Lin. Do llms recognize your preferences? evaluating personalized preference following in llms. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025*. OpenReview.net, 2025d. https://openreview.net/forum?id=QWunLKbBGF.

Boyuan Zheng, Michael Y. Fatemi, Xiaolong Jin, Zora Zhiruo Wang, Apurva Gandhi, Yueqi Song, Yu Gu, Jayanth Srinivasa, Gaowen Liu, Graham Neubig, and Yu Su. Skillweaver: Web agents can self-improve by discovering and honing skills. *CoRR*, abs/2504.07079, 2025a. doi: 10.48550/ARXIV.2504.07079. https://doi.org/10.48550/arXiv.2504.07079.

Junhao Zheng, Xidi Cai, Qiuke Li, Duzhen Zhang, ZhongZhi Li, Yingying Zhang, Le Song, and Qianli Ma. Lifelongagentbench: Evaluating llm agents as lifelong learners. *arXiv preprint arXiv:2505.11942*, 2025b.

Longtao Zheng, Rundong Wang, Xinrun Wang, and Bo An. Synapse: Trajectory-as-Exemplar Prompting with Memory for Computer Control. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024a. https://openreview.net/forum?id=Pc8AU1aF5e.

Yicong Zheng, Kevin L. McKee, Thomas Miconi, Zacharie Bugaud, Mick van Gelderen, and Jed McCaleb. Goal-directed search outperforms goal-agnostic memory compression in long-context memory tasks, 2025c. https://arxiv.org/abs/2511.21726.

Yuanhang Zheng, Peng Li, Wei Liu, Yang Liu, Jian Luan, and Bin Wang. ToolRerank: Adaptive and Hierarchy-Aware Reranking for Tool Retrieval. In Nicoletta Calzolari, Min-Yen Kan, Véronique Hoste, Alessandro Lenci, Sakriani Sakti, and Nianwen Xue, editors, *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation, LREC/COLING 2024, 20-25 May, 2024, Torino, Italy*, pages 16263–16273. ELRA and ICCL, 2024b. https://aclanthology.org/2024.lrec-main.1413.

Siru Zhong, Weilin Ruan, Ming Jin, Huan Li, Qingsong Wen, and Yuxuan Liang. Time-vlm: Exploring multimodal vision-language models for augmented time series forecasting. *CoRR*, abs/2502.04395, 2025. doi: 10.48550/ARXIV.2502.04395. https://doi.org/10.48550/arXiv.2502.04395.

Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. Memorybank: Enhancing large language models with long-term memory. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 19724–19731, 2024.

Huichi Zhou, Yihang Chen, Siyuan Guo, Xue Yan, Kin Hei Lee, Zihan Wang, Ka Yiu Lee, Guchun Zhang, Kun Shao, Linyi Yang, and Jun Wang. Memento: Fine-tuning LLM agents without fine-tuning llms. *CoRR*, abs/2508.16153, 2025a. doi: 10.48550/ARXIV.2508.16153. https://doi.org/10.48550/arXiv.2508.16153.

Jinfeng Zhou, Zhuang Chen, Dazhen Wan, Bosi Wen, Yi Song, Jifan Yu, Yongkang Huang, Pei Ke, Guanqun Bi, Libiao Peng, Jiaming Yang, Xiyao Xiao, Sahand Sabour, Xiaohan Zhang, Wenjing Hou, Yijia Zhang, Yuxiao Dong, Hongning Wang, Jie Tang, and Minlie Huang. Characterglm: Customizing social characters with large language models. In Franck Dernoncourt, Daniel Preotiuc-Pietro, and Anastasia Shimorina, editors, *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: EMNLP 2024 - Industry Track, Miami, Florida, USA, November 12-16, 2024*, pages 1457–1476. Association for Computational Linguistics, 2024a. doi: 10.18653/V1/2024.EMNLP-INDUSTRY.107. https://doi.org/10.18653/v1/2024.emnlp-industry.107.

Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, et al. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023a.

Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan

Bisk, Daniel Fried, Uri Alon, and Graham Neubig. Webarena: A realistic web environment for building autonomous agents, 2024b. https://arxiv.org/abs/2307.13854.

Wangchunshu Zhou, Yuchen Eleanor Jiang, Peng Cui, Tiannan Wang, Zhenxin Xiao, Yifan Hou, Ryan Cotterell, and Mrinmaya Sachan. Recurrentgpt: Interactive generation of (arbitrarily) long text, 2023b. https://arxiv.org/abs/2305.13304.

Wangchunshu Zhou, Yuchen Eleanor Jiang, Long Li, Jialong Wu, Tiannan Wang, Shi Qiu, Jintian Zhang, Jing Chen, Ruipu Wu, Shuai Wang, Shiding Zhu, Jiyu Chen, Wentao Zhang, Ningyu Zhang, Huajun Chen, Peng Cui, and Mrinmaya Sachan. Agents: An open-source framework for autonomous language agents. *CoRR*, abs/2309.07870, 2023c. doi: 10.48550/ARXIV.2309.07870. https://doi.org/10.48550/arXiv.2309.07870.

Wangchunshu Zhou, Yixin Ou, Shengwei Ding, Long Li, Jialong Wu, Tiannan Wang, Jiamin Chen, Shuai Wang, Xiaohua Xu, Ningyu Zhang, Huajun Chen, and Yuchen Eleanor Jiang. Symbolic learning enables self-evolving agents. *CoRR*, abs/2406.18532, 2024c. doi: 10.48550/ARXIV.2406.18532. https://doi.org/10.48550/arXiv.2406.18532.

Zijian Zhou, Ao Qu, Zhaoxuan Wu, Sunghwan Kim, Alok Prakash, Daniela Rus, Jinhua Zhao, Bryan Kian Hsiang Low, and Paul Pu Liang. MEM1: Learning to synergize memory and reasoning for efficient long-horizon agents. *arXiv preprint arXiv:2506.15841*, 2025b.

Andrew Zhu, Lara Martin, Andrew Head, and Chris Callison-Burch. CALYPSO: LLMs as dungeon masters' assistants. In *Proceedings of the Nineteenth AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, volume 19 of *AIIDE '23*, pages 380–390, Salt Lake City, October 2023. AAAI Press. ISBN 978-1-57735-883-1. doi: 10.1609/aiide.v19i1.27534.

Wazeer Deen Zulfikar, Samantha W. T. Chan, and Pattie Maes. Memoro: Using large language models to realize a concise interface for real-time memory augmentation. In Florian 'Floyd' Mueller, Penny Kyburz, Julie R. Williamson, Corina Sas, Max L. Wilson, Phoebe O. Toups Dugas, and Irina Shklovski, editors, *Proceedings of the CHI Conference on Human Factors in Computing Systems, CHI 2024, Honolulu, HI, USA, May 11-16, 2024*, pages 450:1–450:18. ACM, 2024. doi: 10.1145/3613904.3642450. https://doi.org/10.1145/3613904.3642450.

Jialong Zuo, Yongtai Deng, Lingdong Kong, Jingkang Yang, Rui Jin, Yiwei Zhang, Nong Sang, Liang Pan, Ziwei Liu, and Changxin Gao. Videolucy: Deep memory backtracking for long video understanding, 2025. https://arxiv.org/abs/2510.12422.