

# MEM<sup>2</sup>EVOLVE: TOWARDS SELF-EVOLVING AGENTS VIA CO-EVOLUTIONARY CAPABILITY EXPANSION AND EXPERIENCE DISTILLATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

While large language model–powered agents can self-evolve by accumulating experience or by dynamically creating new assets (i.e., tools or expert agents), existing frameworks typically treat these two evolutionary processes in isolation. This separation overlooks their intrinsic interdependence: the former is inherently bounded by a manually predefined static toolset, while the latter generates new assets from scratch without experiential guidance, leading to limited capability growth and unstable evolution. To address this limitation, we introduce a novel paradigm of co-evolutionary Capability Expansion and Experience Distillation. Guided by this paradigm, we propose the **Mem<sup>2</sup>Evolve**, which integrates two core components: **Experience Memory** and **Asset Memory**. Specifically, Mem<sup>2</sup>Evolve leverages accumulated experience to guide the dynamic creation of assets, thereby expanding the agent’s capability space while simultaneously acquiring new experience to achieve co-evolution. Extensive experiments across 6 task categories and 8 benchmarks demonstrate that Mem<sup>2</sup>Evolve achieves improvement of 18.53% over standard LLMs, 11.80% over agents evolving solely through experience, and 6.46% over those evolving solely through asset creation, establishing it as a substantially more effective and stable self-evolving agent framework<sup>1</sup>.

## 1 INTRODUCTION

Large language model (LLM)–powered agents have achieved remarkable success in a wide range of applications (Yang et al., 2024; Jin et al., 2025; Deng et al., 2025; Liu et al., 2025). Building on these successes, recent research is moving beyond static, task-specific systems toward self-evolving agents that can leverage past experiences and autonomously expand their capabilities (Gao et al., 2025; Fang et al., 2025).

However, current frameworks predominantly treat these evolutionary processes in isolation (Cemri et al., 2025). As illustrated in Figure 1(a), **Experience-centric evolution** (Yuan et al., 2025; Yuksekogonul et al., 2025) enables systems to learn from experience to optimize execution strategies (Ma et al., 2025), refine prompts (Zhang et al., 2025a), or build experience repositories (Ouyang et al., 2025). However, this paradigm limits the system to a fixed set of tools and expert agents, leading capability space remains static and cannot expand beyond the pre-specified library. In contrast, **capability-centric evolution** (Figure 1(b)) enables the system to dynamically create new tools (Wölfllein et al., 2025; Qiu et al., 2025) or spawn new expert agents (Chen et al., 2023; 2024; Zhang et al., 2025b). However, creating new assets from scratch without the guidance of experience prevents the system from utilizing proven strategies and avoiding known pitfalls, leading to non-replicable success and repeated errors.

To address these limitations, inspired by the equilibrium theory (Piaget, 1972), which posits intelligence evolves through the interplay of assimilation (integrating new experiences) and accommodation (adapting internal structures), we introduce a novel paradigm of *co-evolutionary capability expansion and experience distillation* (Figure 1c). In this paradigm, expanding agent capabilities enables it to complete a broader range of tasks, thereby yielding more experiences. These experiences are

<sup>1</sup>Our code will be publicly available.

then distilled to guide subsequent capability expansion, realizing co-evolution of capability and experience.

Guided by this paradigm, we propose **Mem<sup>2</sup>Evolve**, an agentic framework that coordinates the evolution of capabilities and experiences through a core dual-memory mechanism comprising Asset Memory and Experience Memory. Specifically, Asset Memory serves as a persistent and extensible repository of the agent’s capabilities, organizing expert agents and executable tools. Experience Memory accumulates strategic experience distilled from both successful and failed trajectories to guide future asset creation and task execution. Building upon this dual-memory architecture, Mem<sup>2</sup>Evolve operates through two complementary phases. During *forward inference*, the system follows a “reuse first, create on demand” strategy, leveraging both memories to execute tasks. When a task exceeds the agent’s current capability boundary, the system dynamically creates new assets guided by experience to expand its capabilities. Upon task completion, *backward evolution* retains high-quality newly created assets into Asset Memory and distills transferable lessons into Experience Memory. This forward-backward loop enables the co-evolution of capabilities and experiences.

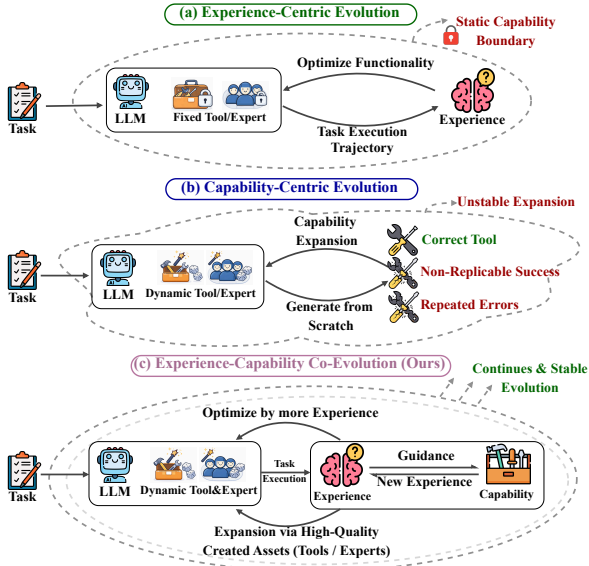


Figure 1: **Paradigms of Self-Evolving Agents:** (a) Experience-centric evolution, (b) Capability-centric evolution, and (c) Our co-evolutionary framework that jointly expands capabilities and distills experience.

To validate the effectiveness of Mem<sup>2</sup>Evolve, we conduct extensive experiments across 6 tasks and 8 benchmarks, covering general assistant (Mialon et al., 2024), multi-hop question answering (Yang et al., 2018), mathematical reasoning, embodied task (Shridhar et al., 2020), planning (Xie et al., 2024), and web interaction (Yao et al., 2022a). Beyond achieving superior overall performance against capability- and experience-centric baselines, Mem<sup>2</sup>Evolve demonstrates robust adaptability, enabling sustained evolution in *single-task* and effective memory reuse in *cross-task* settings.

Our contributions are summarized as follows:

- To the best of our knowledge, we are the first to propose the co-evolutionary agent paradigm that couples dynamic capability expansion with experience distillation.
- Guided by this paradigm, we introduce Mem<sup>2</sup>Evolve, a dual-memory framework that coordinates Asset Memory for dynamic capability expansion and Experience Memory for strategic experience distillation. Through a forward inference and backward evolution loop, Mem<sup>2</sup>Evolve continuously leverages and expands both memories, driving capability–experience co-evolution.
- Extensive experiments show that Mem<sup>2</sup>Evolve consistently outperforms both capability-centric and experience-centric baselines. Moreover, it exhibits strong adaptability, supporting sustained self-evolution in single-task settings and effective memory reuse for cross-task generalization.

## 2 RELATED WORK

**Experience-Centric Evolving.** Recent research on self-evolving agents predominantly focuses on optimizing systems by leveraging experience accumulated from past tasks (Yuan et al., 2025; Li et al., 2023; Ma et al., 2025). For instance, DyLAN (Liu et al., 2023) and DSPy (Khattab et al., 2023) dynamically select agent teams from a predefined pool by aligning past experience with current task requirements. Similarly, Aflow (Zhang et al., 2025a) and AgentSquare (Shang et al., 2025) modularize

Table 1: **Comparison of self-evolving agent frameworks.** **Optimization** indicates whether experience is used to optimize the agent (e.g., prompts). **Persistence** denotes whether experiences are persistently stored for future reuse. **Source:** 🗃️ agent task execution trajectory, 🛠️ tool creation process. **Tool Crea.** and **Agent Crea.** indicate whether the framework supports creation of tools and expert agents, respectively. **Tool/Agent** denotes whether the toolset and expert agents are static or dynamic. **Crea. Grounding** indicates the knowledge sources used for asset creation, 🌐 parametric knowledge, 🌐 web search information, 🧠 experience. **Exp.-Guided Creation** indicates whether new assets are created under the guidance of past experience. Details in the Appendix A.1 and A.2.

Framework	Experience Distillation			Capability Expansion				Exp.-Guided Creation
	Optimization	Persistence	Source	Tool Crea.	Agent Crea.	Tool/Agent	Crea. Grounding	
DSPy (Khatab et al., 2023)	✓	✗	🗃️	✗	✗	Static	–	✗
DyLAN (Liu et al., 2023)	✓	✗	🗃️	✗	✗	Static	–	✗
ReasoningBank (Ouyang et al., 2025)	✗	✓	🗃️	✗	✗	Static	–	✗
AFlow (Zhang et al., 2025a)	✓	✗	🗃️	✗	✗	Static	–	✗
AgentSquare (Shang et al., 2025)	✓	✗	🗃️	✗	✗	Static	–	✗
Agentic Neural Networks (Ma et al., 2025)	✓	✗	🗃️	✗	✗	Static	–	✗
AgentVerse (Chen et al., 2023)	✓	✗	–	✗	✓	Dynamic	🌐	✗
AutoAgents (Chen et al., 2024)	✗	✗	–	✗	✓	Dynamic	🌐	✗
SwarmAgentic (Zhang et al., 2025b)	✓	✗	–	✗	✓	Dynamic	🌐	✗
Alita (Qiu et al., 2025)	✗	✗	–	✓	✗	Dynamic	🌐+🌐	✗
ToolMaker (Wölflin et al., 2025)	✗	✗	–	✓	✗	Dynamic	🌐+🌐	✗
<b>Mem<sup>2</sup>Evolve (Ours)</b>	✓	✓	🗃️+🛠️	✓	✓	<b>Dynamic</b>	🌐+🌐+🧠	✓

agents and employ search algorithms to optimize module compositions. ReasoningBank (Ouyang et al., 2025) summarizes successful and failed experiences to enhance performance on new tasks. However, as shown in Table 1, these frameworks are confined to a fixed set of tools and agents, resulting in a static capability space. Consequently, they cannot extend their boundaries to handle tasks beyond the predefined asset. In contrast, Mem<sup>2</sup>Evolve dynamically creates high-quality agents and tools, enabling it to transcend these pre-existing capability limits.

**Capability-Centric Evolving.** In parallel to experience-centric evolving, capability-centric frameworks focus on expanding the boundaries of agentic systems by dynamically generating tools or agents, thereby reducing dependence on manual design (Cai et al., 2024; Song et al., 2024). AgentVerse (Chen et al., 2023) and AutoAgents (Chen et al., 2024) generate expert agents tailored to specific task dynamics, extending the system’s execution capabilities. ToolMaker (Wölflin et al., 2025) and Alita (Qiu et al., 2025) dynamically create tools to handle videos, documents, and complex mathematical simulations (Feng et al., 2025). However, creating these new assets from scratch without the guidance of experience prevents these systems from leveraging proven strategies or avoiding known pitfalls. This isolation inevitably leads to non-replicable successes and recurring errors. In contrast, Mem<sup>2</sup>Evolve couples capability expansion with experience distillation, realizing a co-evolution that past insights guide asset creation and new capabilities yield richer experiences.

### 3 MEM<sup>2</sup>EVOLVE

We present **Mem<sup>2</sup>Evolve**, a novel self-evolving agent framework that coordinates capability expansion and experience distillation. As illustrated in Figure 2, Mem<sup>2</sup>Evolve is built upon a *Dual-Memory Mechanism: Asset Memory* for dynamic capability expansion and *Experience Memory* for strategic experience distillation (§3.1). Built on this dual-memory foundation, Mem<sup>2</sup>Evolve operates in a two-phase task loop: *forward inference* and *backward evolution*. During *forward inference* (§3.2), the agent leverages both memories to execute tasks while dynamically creating new assets to expand its capabilities boundary. Upon task completion, the *backward evolution* process (§3.3) retains high-quality assets and distills lessons from execution trajectories, enabling continuous self-evolution.

#### 3.1 DUAL-MEMORY MECHANISM

We organize the memory into two distinct components: the **Asset Memory**  $\mathcal{M}_A$ , which stores the expert agents and tools, and the **Experience Memory**  $\mathcal{M}_E$ , which accumulates lessons distilled from past successes and failures to guide future actions.

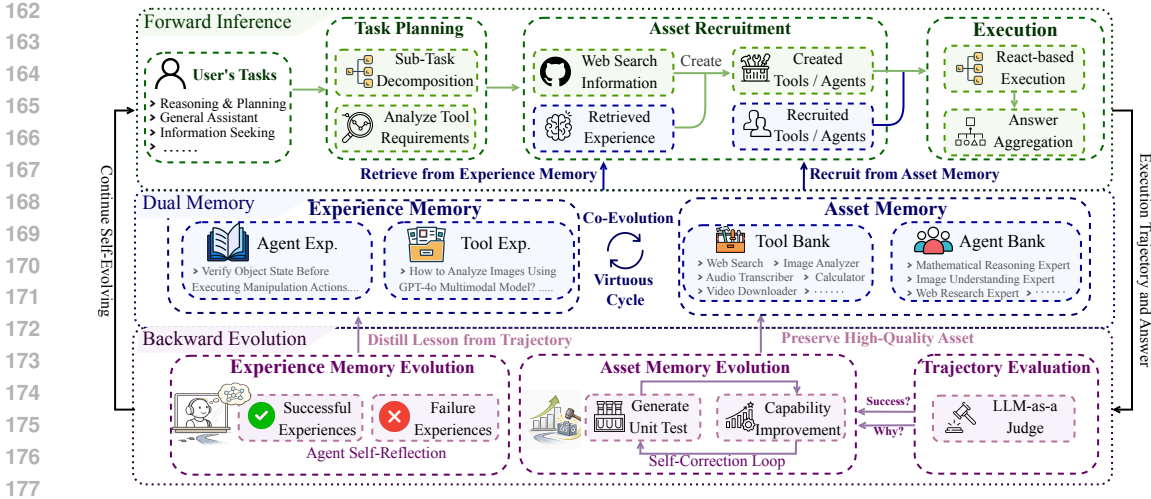


Figure 2: **Overview of Mem<sup>2</sup>Evolve**, a self-evolving agent framework built on a **Dual-Memory** mechanism. The evolution proceeds in two phases. During **Forward Inference**, the agent recruits tools and expert agents from *Asset Memory* to execute the current task. When the task exceeds its current capability boundary, *Experience Memory* is leveraged to guide the stable creation of new assets on demand. During **Backward Evolution**, newly validated assets are preserved in *Asset Memory* to achieve persistent capability expansion, while strategic insights distilled from execution trajectories are accumulated into *Experience Memory*. This forward-backward loop enables the co-evolution of capabilities and experience, forming a stable self-evolving cycle.

### 3.1.1 ASSET MEMORY

To support capability expansion at both the strategic level (through expert agents) and the operational level (through tools), *Asset Memory* serves as a repository of reusable, execution-ready capabilities:

$$\mathcal{M}_A = \mathcal{B}_{agt} \cup \mathcal{B}_{tool}, \quad (1)$$

where  $\mathcal{B}_{agt}$  is the *Agent Bank* containing expert agents, and  $\mathcal{B}_{tool}$  is the *Tool Bank* that stores executable tools.

**Agent Bank.** Building on prior work (Chen et al., 2024) and Anthropic’s Agent Skills<sup>2</sup>, we distill a compact *agent specification* tailored to Mem<sup>2</sup>Evolve. As exemplified in Figure 5, each entry  $m_{agt} \in \mathcal{B}_{agt}$  is defined as:

$$m_{agt} = \langle \rho, \epsilon, \sigma, \mathbb{T}_{avail} \rangle, \quad (2)$$

where  $\rho$  is the *role* specifying the agent’s identity,  $\epsilon$  describes the agent’s *expertise* and domain knowledge,  $\sigma$  denotes *suggestions* that guide the agent’s behavior strategies, and  $\mathbb{T}_{avail} \subseteq \mathcal{B}_{tool}$  specifies the set of available tools.

**Tool Bank.** To ensure seamless integration with diverse LLM backbones, *Tool Bank* maintains executable tools stored in compliance with the Model Context Protocol (MCP)<sup>3</sup>, example in Code 1. Each entry  $m_{tool} \in \mathcal{B}_{tool}$  is defined as:

$$m_{tool} = \langle n, d_{func}, c_{impl}, \omega_{doc} \rangle, \quad (3)$$

where  $n$  is the *tool name*,  $d_{func}$  provides a *functional description*,  $c_{impl}$  contains the *implementation code*, and  $\omega_{doc}$  specifies input/output documentation.

### 3.1.2 EXPERIENCE MEMORY

To enable Mem<sup>2</sup>Evolve to replicate proven strategies and circumvent previously encountered pitfalls, *Experience Memory* accumulates insights derived from past successes and failures, guiding future

<sup>2</sup><https://github.com/anthropics/skills>

<sup>3</sup><https://www.anthropic.com/news/model-context-protocol>

task execution and asset creation (Ouyang et al., 2025). We define  $\mathcal{M}_E = \mathcal{E}_{agt} \cup \mathcal{E}_{tool}$ , with each *Memory Item*  $e \in \mathcal{M}_E$  structured as:

$$e = \langle h_{title}, d_{desc}, \mathcal{U}_{case}, \kappa_{content} \rangle, \quad (4)$$

where  $h_{title}$  is the title,  $d_{desc}$  describes the context,  $\mathcal{U}_{case}$  lists applicable use cases, and  $\kappa_{content}$  stores the core distilled knowledge, encompassing both agent experience and tool experience:

**Agent Experience.**  $\kappa_{content}$  contains strategic insights derived from trajectory reflections, guiding specific expert agents in handling complex tasks.

**Tool Experience.**  $\kappa_{content}$  contains implementation guidelines distilled from the tool creation and debugging process, and an example in Figure 8.

## 3.2 FORWARD INFERENCE

To balance the utilization of accumulated expertise with the acquisition of new capabilities, the forward inference follows a strategy of "Reuse first, Create on demand". We formalize it into three phases: (1) *task planning*, (2) *asset recruitment*, and (3) *execution*.

### 3.2.1 TASK PLANNING

Initially, the LLM  $\pi_\theta$  acts as a planner to decompose the task  $q_t$  into a sequence of sub-tasks  $\mathcal{S} = \{s_1, s_2, \dots, s_k\}$ , with the prompt in Appendix C.1. This decomposition ensures that complex problems are broken down into solvable units with clear resource definitions.

### 3.2.2 ASSET RECRUITMENT

For each sub-task  $s_i$ , the system prepares the required assets via the Recruitment Function  $\Gamma(s_i)$ :

$$\Gamma(s_i) = \begin{cases} m^* & \text{sim}(s_i, \mathcal{M}_A) \geq \delta \\ \text{Create}(s_i | \mathcal{M}_E, \text{Web}) & \text{otherwise} \end{cases} \quad (5)$$

where  $\text{sim}(s_i, \mathcal{M}_A)$  measures the similarity between the sub-task and the asset stored in Assets Memory, and  $\delta$  is a confidence threshold. This mechanism determines whether the sub-task lies beyond the agent’s current capability boundary. Depending on the output of  $\Gamma(s_i)$ , the process branches into two paths:

**Recruitment.** If a high-similarity match exists, the system directly reuses  $m^* \in \mathcal{M}_A$ . For agents, we select the top-1 candidate surpassing  $\delta$  to entrust the sub-task to the most specialized expert. Conversely, for tools, we retrieve the top- $k$  matches to ensure comprehensive utility while mitigating context overhead from excessive documentation.

**Creation.** Conversely, for missing capabilities, *Tool Creation* employs experience-augmented generation, conditioning on the sub-task  $s_i$ , web search results, and relevant experiences  $e$  from  $\mathcal{E}_{tool}$ :

$$m_{tool}^{new} \sim \pi_\theta(s_i | \text{Retrieve}(s_i, \mathcal{E}_{tool}), \text{Web}(s_i)). \quad (6)$$

Similarly, *Agent Creation* synthesizes a new expert by prompting  $\pi_\theta$  with task requirements derived from  $s_i$ , and details in Appendix A.3.

### 3.2.3 EXECUTION

Each sub-task  $s_i$  is assigned to its recruited agent  $m_{agt}^i$ , augmented with experiences  $e$  retrieved from  $\mathcal{E}_{agt}$  for role-specific guidance. The agent then executes using available tools  $\mathbb{T}_{avail}^i$  within the ReAct framework (Yao et al., 2022b), alternating among think, action, and observation steps. Finally, system aggregates all results  $\{r_1, \dots, r_k\}$  to produce the final answer  $a_t$ .

## 3.3 BACKWARD EVOLUTION

Upon task completion, the backward evolution aims to preserve high-quality assets for future reuse and distill transferable lessons from execution trajectories. We formalize it into: (1) *trajectory evaluation*, (2) *asset memory evolution*, and (3) *experience memory evolution*.

### 3.3.1 TRAJECTORY EVALUATION

The evaluation stage provides the foundation for all subsequent memory updates. We employ an LLM-as-a-Judge (Li et al., 2025a; Ouyang et al., 2025) to assess execution quality.<sup>4</sup> Given the task  $q_t$ , execution trajectory  $\tau_t$ , and answer  $a_t$ , the Judge produces:

$$r_t, c_t = \text{Judge}(q_t, \tau_t, a_t), \quad (7)$$

where  $r_t \in \{0, 1\}$  indicates success or failure, and  $c_t$  provides critique comments identifying specific strengths and weaknesses.

### 3.3.2 ASSET MEMORY EVOLUTION

This phase determines which newly created assets should be preserved and refined before entering  $\mathcal{M}_A$ . Since a correct answer does not guarantee robust underlying assets, we adopt a *Self-Correction Loop* guided by  $r_t$  and  $c_t$ .

For each asset  $m_{\text{new}} \in \mathcal{A}_t^{\text{new}}$ , where  $\mathcal{A}_t^{\text{new}}$  denotes the set of newly created assets, we derive a finalized version  $m_{\text{final}}$  as:

$$m_{\text{final}} = \begin{cases} m_{\text{new}} & \text{if } r_t = 1 \wedge \\ & \text{Valid}(m_{\text{new}}, c_t) \\ \text{Improve}(m_{\text{new}}, c_t) & \text{otherwise} \end{cases} \quad (8)$$

Where  $\text{Valid}(m_{\text{new}}, c_t)$  verifies asset reliability by having the LLM synthesize test cases from the critique  $c_t$  and executing  $m_{\text{new}}$  against them. An asset passes validation only if it clears all tests.

If validation fails,  $\text{Improve}(m_{\text{new}}, c_t)$  triggers a *Self-Correction Loop*: revise the asset based on  $c_t$  and test failures, then regenerate tests until validation passes. Once validated:

$$\mathcal{M}_A \leftarrow \mathcal{M}_A \cup \{m_{\text{final}}\}. \quad (9)$$

### 3.3.3 EXPERIENCE MEMORY EVOLUTION

Mem<sup>2</sup>Evolve distills trajectory-level insights into the  $\mathcal{M}_E$  to guide future task execution and asset creation. After each task, the system reflects on the trajectory  $\tau_t$  and  $(r_t, c_t)$  to extract *Memory Items*:

$$e_{\text{new}} = \text{Reflection}(\tau_t, r_t, c_t), \quad (10)$$

The Reflection function captures insights from both successful and failed executions.

**Success Generalization.** When  $r_t = 1$ , Mem<sup>2</sup>Evolve abstracts high-level guidelines from the successful trajectory. For agents,  $\kappa_{\text{content}}$  records strategic advice and coordination patterns for specific roles; for tools, it captures effective implementation patterns and usage recipes.

**Failure Diagnosis.** When  $r_t = 0$  or  $c_t$  indicates substantial debugging, Reflection focuses on failure modes. The resulting  $e_{\text{new}}$  encodes anti-patterns and failure-fix pairs to prevent similar errors. Detailed prompt in Appendix C.8 and C.9.

Finally, the distilled experience items are merged into the Experience Memory:

$$\mathcal{M}_E \leftarrow \mathcal{M}_E \cup \{e_{\text{new}}\}. \quad (11)$$

## 4 EXPERIMENTS

### 4.1 EXPERIMENT SETTING

**Baselines.** Following prior research (Qiu et al., 2025; Zhang et al., 2025b), we compare Mem<sup>2</sup>Evolve against three categories of baselines: (1) **Naive LLMs**, including Direct prompting, CoT (Wei et al., 2022), ReAct (Yao et al., 2022b), and OpenAI’s DeepResearch (OpenAI, 2025),

<sup>4</sup>We assume ground-truth labels are inaccessible during backward evolution to simulate real world. When available, such supervision can further enhance evolution effectiveness.

Table 2: **Main results across 6 tasks and 8 benchmarks**, reported as Pass@1 for each benchmark. The best results are highlighted in **bold**, and the second-best results are underlined. †Results are from the original paper.

Method	GAIA				Embodied	Multi-Hop QA		Math		Planning	Web Interaction	Avg.
	L1	L2	L3	Total	ALFWorld	HotpotQA	2Wiki	AIME24	AIME25	TravelPlanner	WebShop	
<i>Naive-Large Language Model</i>												
GPT-5-Chat (Direct)	16.98	12.79	7.69	12.49	83.58	50.40	<u>81.80</u>	60.00	46.67	38.68	22.31	49.49
GPT-5-Chat (CoT)	24.53	17.44	11.54	17.84	83.58	47.40	74.40	66.67	56.67	39.51	27.49	51.71
GPT-5-Chat (ReAct)	26.42	17.44	11.54	18.47	86.87	41.40	48.40	66.67	60.00	39.13	25.10	48.27
OpenAI-DeepResearch†	74.29	69.06	<u>47.60</u>	67.36	—	—	—	—	—	—	—	—
<i>Experience-Centric Evolving</i>												
DyLAN	24.53	19.78	11.54	18.62	91.20	52.00	65.00	46.67	43.33	43.15	36.40	49.55
EvoAgent	22.64	19.78	11.54	17.99	92.50	54.40	75.00	66.67	43.33	49.20	37.80	54.61
AFLOW	26.42	17.44	15.38	19.75	<u>93.40</u>	<b>60.80</b>	72.40	66.67	63.33	53.24	<u>37.90</u>	58.44
DSPy	30.19	15.12	11.54	18.95	92.80	55.60	76.40	66.67	50.00	44.90	35.50	55.10
<i>Capacity-Centric Evolving</i>												
Alita	<u>81.13</u>	<u>75.58</u>	46.15	<u>72.73</u>	86.13	<u>58.80</u>	77.40	<u>70.00</u>	<u>66.67</u>	48.32	30.21	<u>63.78</u>
AgentVerse	30.19	16.28	19.23	21.90	88.32	38.60	74.60	60.00	50.00	47.25	32.53	51.65
AutoAgents	35.85	24.42	19.23	26.50	87.92	54.20	73.80	40.00	36.67	43.52	31.40	49.25
SwarmAgentic	28.30	18.60	13.46	20.40	88.79	56.00	80.00	46.67	40.00	<u>59.14</u>	34.12	53.14
<i>Ours</i>												
Mem <sup>2</sup> Evolve	<b>88.68</b>	<b>82.56</b>	<b>57.69</b>	<b>76.31</b>	<b>94.31</b>	<b>60.80</b>	<b>82.00</b>	<b>76.70</b>	<b>73.33</b>	<b>59.25</b>	<b>39.20</b>	<b>70.24</b>

(2) **Experience-Centric frameworks** such as DyLAN (Liu et al., 2023), EvoAgent (Yuan et al., 2025), AFLOW (Zhang et al., 2025a), and DSPy (Khatab et al., 2023), and (3) **Capacity-Centric frameworks**, spanning tool-generative methods Alita (Qiu et al., 2025), ToolMaker (Wölflein et al., 2025) and agent-generative approaches AgentVerse (Chen et al., 2023), AutoAgents (Chen et al., 2024), SwarmAgentic (Zhang et al., 2025b). More details in the Appendix B.1.

**Benchmarks.** Following Li et al. (2025b), we evaluate the agent’s capabilities across 8 benchmarks in 6 distinct tasks. These include **GAIA** (Mialon et al., 2024) for general assistant, **ALF-World** (Shridhar et al., 2020) and **WebShop** (Yao et al., 2022a) for embodied and web interaction, and **TravelPlanner** (Xie et al., 2024) for planning. We also include **HotpotQA** (Yang et al., 2018) and **2WikiMultihopQA** (Ho et al., 2020) for multi-hop QA, plus **AIME 24/25** for mathematical reasoning. Details are in Appendix B.2.

**Implement Details.** For all baselines, we utilize GPT-5-chat<sup>5</sup> as the LLM backbone. The web search tool incorporates the Serper search engine<sup>6</sup> and the Crawl4AI (UncleCode, 2024) parsing framework, and code execution is managed via the SandboxFusion environment (Bytedance-Seed-Foundation-Code-Team et al., 2025).

## 4.2 MAIN RESULTS

Table 2 presents the comparative results of different frameworks, and the following conclusions are derived based on these results.

**Capability-experience co-evolution achieves the strongest general agent.** Mem<sup>2</sup>Evolve achieves the best overall performance among all evaluated frameworks, demonstrating the effectiveness of jointly evolving capabilities and experience. Under the same GPT-5-chat as all baselines, Mem<sup>2</sup>Evolve attains an average Pass@1 of 70.24% across all benchmarks, outperforming the strongest capability-centric baseline Alita by 6.46%, experience-centric baseline Aflow by 11.80%, and naive-llm by up to 18.53%. These consistent improvements confirm that capability–experience co-evolution yields a substantially more powerful general agent than either paradigm.

**Breaking the Capability Boundaries of Static Agents.** When initialized with only a Web Search, purely experience-centric methods yield marginal improvements over the base LLM. On GAIA,

<sup>5</sup><https://openai.com/index/introducing-gpt-5/>

<sup>6</sup><https://serpapi.com/>

experience-centric baselines with a fixed toolset improve Pass@1 by at most 1.28%; on AIME, AFLOW achieves only +3.33% on AIME25 and no improvement on AIME24. In contrast, Mem<sup>2</sup>Evolve, starting from the same minimal configuration but capable of evolving new tools and expert agents, achieves +57.84% on GAIA and +10.03%/+13.33% on AIME24/AIME25, respectively. These substantial gains demonstrate that Mem<sup>2</sup>Evolve effectively extends the capability boundary of the base LLM.

**Experience Memory Enhances Capability Expansion.** Incorporating Experience Memory further enhances the effectiveness of capability expansion. Under matched conditions, Mem<sup>2</sup>Evolve outperforms the capability-centric baseline Alita by 6.46% in average Pass@1. This improvement suggests that Experience Memory refines and stabilizes the utilization of newly evolved tools and agents, enabling capability expansion to translate more reliably into downstream performance gains.

## 5 ANALYSIS

In this section, we conduct a comprehensive analysis to answer the following research questions **RQ1: What role does each module play in Mem<sup>2</sup>Evolve?** (§5.1) **RQ2: How does experience guide asset generation?** (§5.2) **RQ3: How does Mem<sup>2</sup>Evolve self-evolve in single task?** (§5.3) **RQ4: How does Mem<sup>2</sup>Evolve self-evolve across tasks?** (§5.4) **RQ5: What is being evolved during self-evolution?** (§B.4) **RQ6: How does Mem<sup>2</sup>Evolve behave in case studies?** (§D)

### 5.1 RQ1: ABLATION STUDY

To verify the effectiveness of each module in Mem<sup>2</sup>Evolve, we conducted an ablation study on Asset Creation and Experience Distillation. As shown in Table 3, Mem<sup>2</sup>Evolve consistently outperforms all variants, validating the necessity of the proposed Dual-Memory mechanism. Specifically, *w/o Tool Creation* causes the largest performance drop of 10.28%, highlighting that dynamically expanding the toolset is crucial for handling complex tasks, while *w/o Expert Agent Creation* still leads to a 1.72% decline because all tasks are forced onto a single general-purpose agent rather than expert agents. Moreover, removing *Agent Memory* causes a 4.73% performance drop, and removing *Tool Memory* causes a 3.13% performance drop, as this prevents the system from leveraging validated successes and past failures during both tool creation and task execution, making it difficult to reliably reproduce effective behaviors and avoid known mistakes, thereby degrading overall performance.

Table 3: **Ablation study of Mem<sup>2</sup>Evolve.** Full results are provided in Appendix 6.

Framework	Avg. Pass@1	$\Delta_{rel}^{\%}$
Mem <sup>2</sup> Evolve	70.24	-
<i>w/o Asset Creation</i>		
w/o Tool Creation	59.96	↓ 10.28
w/o Expert Agent Creation	68.52	↓ 1.72
<i>w/o Experience Distillation</i>		
w/o Tool Memory	67.11	↓ 3.13
w/o Agent Memory	65.51	↓ 4.73

### 5.2 RQ2: EXPERIENCE-GUIDED ASSET CREATION

Table 4: **Experience-Guided Asset Creation.** Experience guidance significantly enhances efficiency.

Benchmark	w/o Exp.-Guide	w/ Exp.-Guide	$\Delta_{rel}^{\%}$
<i>First-Pass Validity</i> (↑)			
GAIA	32.7%	51.0% (+18.3%)	↑ 56.0
AIME24	64.9%	83.8% (+18.9%)	↑ 29.1
AIME25	61.8%	82.4% (+20.6%)	↑ 33.3
Avg.	53.1%	72.4% (+19.3%)	↑ 36.3
<i>Avg. Improve Iter.</i> (↓)			
GAIA	1.45	0.94 (-0.51)	↓ 35.2
AIME24	0.76	0.24 (-0.52)	↓ 68.4
AIME25	0.82	0.26 (-0.56)	↓ 68.3
Avg.	1.01	0.48 (-0.53)	↓ 52.5

debugging iterations by nearly 68% and increases first-pass validity to over 82%, indicating more accurate tool generation at the initial attempt. In contrast, in the more complex GAIA, experience

In Section 5.1, we show that incorporating *Tool Memory* leads to consistent performance improvements. This section further investigates its impact on benchmarks that require extensive tool creation. We evaluate experience guidance using: (1) *First-Pass Validity*, which measures whether the initially generated tool satisfies the verification function  $\text{Valid}(m_{\text{new}}, c_t)$  in Equation 8, and (2) *Avg. Improve Iter.*, defined as the average steps of  $\text{Improve}(m_{\text{new}}, c_t)$  during the self-correction loop. As shown in Table 4, experience guidance substantially improves the reliability and efficiency of tool creation. For AIME24/25, where the agent already demonstrates strong performance, experience guidance reduces the average number of

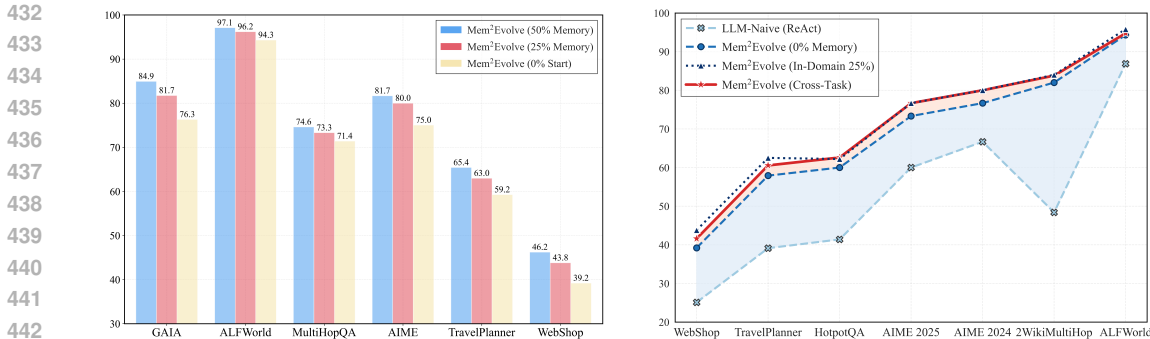


Figure 3: **Self-evolving performance analysis.** (Left) RQ3: Single-task self-evolving: initializing with prior memory consistently improves performance by leveraging accumulated experience. (Right) RQ4: Cross-task self-evolving: heterogeneous memory from GAIA enables stable evolution on target tasks without negative transfer.

guidance improves first-pass validity by 56.0% relative to the w/o Exp-Guide, suggesting that experience effectively constrains tool generation toward feasible solutions. Results and case in Figure 9 show that experience guidance significantly enhances the stability of tool generation, ensuring a more robust evolutionary trajectory for the agent.

### 5.3 RQ3: SINGLE TASK SELF-EVOLVING

In Table 2, we evaluate the performance of Mem<sup>2</sup>Evolve across multiple benchmarks, where each run starts without any pre-existing memory except for access to the web search tool. In this section, we further analyze the effect of introducing initial memory within the same task. Specifically, we construct initial memory using a subset of data from each benchmark and use it to initialize the system, after which evaluation is conducted on the remaining test set. As shown in Figure 3(Left), introducing initial memory consistently improves performance across all benchmarks compared to the setting without initial memory. Most of the performance gains are achieved with a relatively small amount of initial memory, while further enlarging the memory yields diminishing incremental improvements. This pattern suggests that memory accumulated within the same task is effective in enhancing agent performance, with early-stage memory capturing a large fraction of broadly applicable assets and high-utility experience.

### 5.4 RQ4: CROSS TASKS SELF-EVOLVING

To evaluate the generalization capability of Mem<sup>2</sup>Evolve in a cross-task setting, we initialize the agent with heterogeneous memory accumulated from GAIA and evaluate its performance across 7 target benchmarks. As illustrated in Figure 3(Right), cross-task memory initialization consistently improves performance compared to the setting without initial memory, and achieves results comparable to the 25% single-task initialization. Despite the domain mismatch between source and target tasks, the agent maintains stable evolutionary trajectories without suffering negative transfer. These results suggest that Mem<sup>2</sup>Evolve can reuse heterogeneous memory across tasks without adversely affecting performance. The structured representation of memory components and the retrieval mechanism contribute to this behavior by enabling selective access to task-relevant information.

## 6 CONCLUSION

We introduce Mem<sup>2</sup>Evolve, a self-evolving agent framework that integrates Asset Memory and Experience Memory and enables their coordinated co-evolution. This design allows the agent to expand its capability space while continuously accumulating strategic experience, leading to more stable and sustained performance improvements. Extensive experimental results show that Mem<sup>2</sup>Evolve consistently improves performance in both single-task and cross-task settings. We hope that Mem<sup>2</sup>Evolve provides a practical foundation for building general-purpose, lifelong-learning agents with reduced reliance on human intervention.

## REFERENCES

- 486  
487  
488 Bytedance-Seed-Foundation-Code-Team, :, Yao Cheng, Jianfeng Chen, Jie Chen, Li Chen, Liyu  
489 Chen, Wentao Chen, Zhengyu Chen, Shijie Geng, Aoyan Li, Bo Li, Bowen Li, Linyi Li, Boyi  
490 Liu, Jiaheng Liu, Kaibo Liu, Qi Liu, Shukai Liu, Siyao Liu, Tianyi Liu, Tingkai Liu, Yongfei Liu,  
491 Rui Long, Jing Mai, Guanghan Ning, Z. Y. Peng, Kai Shen, Jiahao Su, Jing Su, Tao Sun, Yifan  
492 Sun, Yunzhe Tao, Guoyin Wang, Siwei Wang, Xuwu Wang, Yite Wang, Zihan Wang, Jinxiang Xia,  
493 Liang Xiang, Xia Xiao, Yongsheng Xiao, Chenguang Xi, Shulin Xin, Jingjing Xu, Shikun Xu,  
494 Hongxia Yang, Jack Yang, Yingxiang Yang, Jianbo Yuan, Jun Zhang, Yufeng Zhang, Yuyu Zhang,  
495 Shen Zheng, He Zhu, and Ming Zhu. Fullstack bench: Evaluating llms as full stack coders, 2025.  
496 URL <https://arxiv.org/abs/2412.00535>. 7
- 497  
498 Tianle Cai, Xuezhi Wang, Tengyu Ma, Xinyun Chen, and Denny Zhou. Large language models as  
499 tool makers. In *The Twelfth International Conference on Learning Representations*, 2024. URL  
500 <https://openreview.net/forum?id=qV83K9d5WB>. 3
- 501  
502 Mert Cemri, Melissa Z. Pan, Shuyi Yang, Lakshya A. Agrawal, Bhavya Chopra, Rishabh Tiwari,  
503 Kurt Keutzer, Aditya Parameswaran, Dan Klein, Kannan Ramchandran, Matei Zaharia, Joseph E.  
504 Gonzalez, and Ion Stoica. Why do multi-agent llm systems fail?, 2025. URL <https://arxiv.org/abs/2503.13657>. 1
- 505  
506 Guangyao Chen, Siwei Dong, Yu Shu, Ge Zhang, Jaward Sesay, Börje Karlsson, Jie Fu, and Yemin  
507 Shi. Autoagents: a framework for automatic agent generation. In *Proceedings of the Thirty-Third  
508 International Joint Conference on Artificial Intelligence*, pp. 22–30, 2024. 1, 3, 4, 7, 18, 23
- 509  
510 Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi  
511 Lu, Yi-Hsin Hung, Chen Qian, et al. Agentverse: Facilitating multi-agent collaboration and explor-  
512 ing emergent behaviors. In *The Twelfth International Conference on Learning Representations*,  
513 2023. 1, 3, 7, 17, 23
- 514  
515 Bin Deng, Yizhe Feng, Zeming Liu, Qing Wei, Xiangrong Zhu, Shuai Chen, Yuanfang Guo, and  
516 Yunhong Wang. RETAIL: Towards real-world travel planning for large language models. In  
517 *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pp.  
518 14881–14913, Suzhou, China, November 2025. Association for Computational Linguistics. ISBN  
519 979-8-89176-332-6. doi: 10.18653/v1/2025.emnlp-main.752. URL <https://aclanthology.org/2025.emnlp-main.752/>. 1
- 520  
521 Jinyuan Fang, Yanwen Peng, Xi Zhang, Yingxu Wang, Xinhao Yi, Guibin Zhang, Yi Xu, Bin Wu,  
522 Siwei Liu, Zihao Li, Zhaochun Ren, Nikos Aletras, Xi Wang, Han Zhou, and Zaiqiao Meng. A  
523 comprehensive survey of self-evolving ai agents: A new paradigm bridging foundation models and  
524 lifelong agentic systems, 2025. URL <https://arxiv.org/abs/2508.07407>. 1
- 525  
526 Jiazhan Feng, Shijue Huang, Xingwei Qu, Ge Zhang, Yujia Qin, Baoquan Zhong, Chengquan Jiang,  
527 Jinxin Chi, and Wanjun Zhong. Retool: Reinforcement learning for strategic tool use in llms, 2025.  
528 URL <https://arxiv.org/abs/2504.11536>. 3
- 529  
530 Huan-ang Gao, Jiayi Geng, Wenyue Hua, Mengkang Hu, Xinzhe Juan, Hongzhang Liu, Shilong Liu,  
531 Jiahao Qiu, Xuan Qi, Yiran Wu, et al. A survey of self-evolving agents: On path to artificial super  
532 intelligence. *arXiv preprint arXiv:2507.21046*, 2025. 1
- 533  
534 GitHub. Spec kit: Toolkit to help you get started with spec-driven development. <https://github.com/github/spec-kit>, 2025. Accessed: 2025-12-19. 19
- 535  
536 Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. Constructing a multi-  
537 hop QA dataset for comprehensive evaluation of reasoning steps. In Donia Scott, Nuria Bel,  
538 and Chengqing Zong (eds.), *Proceedings of the 28th International Conference on Computa-  
539 tional Linguistics*, pp. 6609–6625, Barcelona, Spain (Online), December 2020. International  
540 Committee on Computational Linguistics. doi: 10.18653/v1/2020.coling-main.580. URL  
541 <https://aclanthology.org/2020.coling-main.580/>. 7, 23
- 542  
543 Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and  
544 Jiawei Han. Search-r1: Training llms to reason and leverage search engines with reinforcement  
545 learning. *arXiv preprint arXiv:2503.09516*, 2025. 1

- 540 Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vard-  
541 hamanan, Saiful Haq, Ashutosh Sharma, Thomas T Joshi, Hanna Moazam, et al. Dspy: Compiling  
542 declarative language model calls into self-improving pipelines. *arXiv preprint arXiv:2310.03714*,  
543 2023. 2, 3, 7, 16, 22
- 544 Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita  
545 Bhattacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, et al. From generation to judgment:  
546 Opportunities and challenges of llm-as-a-judge. In *Proceedings of the 2025 Conference on*  
547 *Empirical Methods in Natural Language Processing*, pp. 2757–2791, 2025a. 6
- 548 Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Com-  
549 municative agents for "mind" exploration of large language model society. *Advances in Neural*  
550 *Information Processing Systems*, 36:51991–52008, 2023. 2
- 551 Xiaoxi Li, Wenxiang Jiao, Jiarui Jin, Guanting Dong, Jiajie Jin, Yinuo Wang, Hao Wang, Yutao Zhu,  
552 Ji-Rong Wen, Yuan Lu, and Zhicheng Dou. Deepagent: A general reasoning agent with scalable  
553 toolsets, 2025b. URL <https://arxiv.org/abs/2510.21618>. 7
- 554 Jingjing Liu, Zeming Liu, Zihao Cheng, Mengliang He, Xiaoming Shi, Yuhang Guo, Xiangrong Zhu,  
555 Yuanfang Guo, Yunhong Wang, and Haifeng Wang. RepoDebug: Repository-level multi-task and  
556 multi-language debugging evaluation of large language models. In *Findings of the Association*  
557 *for Computational Linguistics: EMNLP 2025*, pp. 23784–23813, Suzhou, China, November 2025.  
558 Association for Computational Linguistics. ISBN 979-8-89176-335-7. doi: 10.18653/v1/2025.  
559 findings-emnlp.1294. URL <https://aclanthology.org/2025.findings-emnlp.1294/>. 1
- 560 Zijun Liu, Yanzhe Zhang, Peng Li, Yang Liu, and Diyi Yang. Dynamic llm-agent network: An  
561 llm-agent collaboration framework with agent team optimization. *arXiv preprint arXiv:2310.02170*,  
562 2023. 2, 3, 7, 16, 22
- 563 Xiaowen Ma, Chenyang Lin, Yao Zhang, Volker Tresp, and Yunpu Ma. Agentic neural networks:  
564 Self-evolving multi-agent systems via textual backpropagation. *arXiv preprint arXiv:2506.09046*,  
565 2025. 1, 2, 3, 17
- 566 Grégoire Mialon, Clémentine Fourier, Thomas Wolf, Yann LeCun, and Thomas Scialom. GAIA:  
567 a benchmark for general AI assistants. In *The Twelfth International Conference on Learning*  
568 *Representations*, 2024. URL <https://openreview.net/forum?id=fibxvahvs3>. 2, 7, 23
- 569 OpenAI. Introducing deep research. [https://openai.com/index/  
570 introducing-deep-research/](https://openai.com/index/introducing-deep-research/), 2025. Accessed: 2026-01-02. 6
- 571 Siru Ouyang, Jun Yan, I Hsu, Yanfei Chen, Ke Jiang, Zifeng Wang, Rujun Han, Long T Le, Samira  
572 Daruki, Xiangru Tang, et al. Reasoningbank: Scaling agent self-evolving with reasoning memory.  
573 *arXiv preprint arXiv:2509.25140*, 2025. 1, 3, 5, 6, 16
- 574 Jean Piaget. Development and learning. *Reading in child behavior and development*, pp. 38–46,  
575 1972. 1
- 576 Jiahao Qiu, Xuan Qi, Tongcheng Zhang, Xinzhe Juan, Jiacheng Guo, Yifu Lu, Yimin Wang, Zixin  
577 Yao, Qihan Ren, Xun Jiang, et al. Alita: Generalist agent enabling scalable agentic reasoning with  
578 minimal predefinition and maximal self-evolution. *arXiv preprint arXiv:2505.20286*, 2025. 1, 3, 6,  
579 7, 17, 22
- 580 Yu Shang, Yu Li, Keyu Zhao, Likai Ma, Jiahe Liu, Fengli Xu, and Yong Li. Agentsquare: Automatic  
581 LLM agent search in modular design space. In *The Thirteenth International Conference on*  
582 *Learning Representations*, 2025. URL <https://openreview.net/forum?id=mPdmDYIq7f>. 2, 3,  
583 16
- 584 Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew  
585 Hausknecht. Alfworld: Aligning text and embodied environments for interactive learning. *arXiv*  
586 *preprint arXiv:2010.03768*, 2020. 2, 7, 23
- 587 Linxin Song, Jiale Liu, Jieyu Zhang, Shaokun Zhang, Ao Luo, Shijian Wang, Qingyun Wu, and  
588 Chi Wang. Adaptive in-conversation team building for language model agents. *arXiv preprint*  
589 *arXiv:2405.19425*, 2024. 3

- 594 UncleCode. Crawl4ai: Open-source llm friendly web crawler & scraper. [https://github.com/](https://github.com/unclecode/crawl4ai)  
595 [unclecode/crawl4ai](https://github.com/unclecode/crawl4ai), 2024. 7  
596
- 597 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny  
598 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*  
599 *neural information processing systems*, 35:24824–24837, 2022. 6, 22
- 600 Georg Wölflein, Dyke Ferber, Daniel Truhn, Ognjen Arandjelovic, and Jakob Nikolas Kather.  
601 LLM agents making agent tools. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and  
602 Mohammad Taher Pilehvar (eds.), *Proceedings of the 63rd Annual Meeting of the Association*  
603 *for Computational Linguistics (Volume 1: Long Papers)*, pp. 26092–26130, Vienna, Austria, July  
604 2025. Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/  
605 2025.acl-long.1266. URL <https://aclanthology.org/2025.acl-long.1266/>. 1, 3, 7, 17, 23  
606
- 607 Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and  
608 Yu Su. Travelplanner: a benchmark for real-world planning with language agents. In *Proceedings*  
609 *of the 41st International Conference on Machine Learning*, pp. 54590–54613, 2024. 2, 7, 23
- 610 John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik Narasimhan,  
611 and Ofir Press. Swe-agent: Agent-computer interfaces enable automated software engineering.  
612 *Advances in Neural Information Processing Systems*, 37:50528–50652, 2024. 1
- 613 Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and  
614 Christopher D. Manning. HotpotQA: A dataset for diverse, explainable multi-hop question answer-  
615 ing. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun’ichi Tsujii (eds.), *Proceedings*  
616 *of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 2369–2380,  
617 Brussels, Belgium, October-November 2018. Association for Computational Linguistics. doi:  
618 10.18653/v1/D18-1259. URL <https://aclanthology.org/D18-1259/>. 2, 7, 23
- 619 Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable  
620 real-world web interaction with grounded language agents. *Advances in Neural Information*  
621 *Processing Systems*, 35:20744–20757, 2022a. 2, 7, 23  
622
- 623 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R Narasimhan, and Yuan  
624 Cao. React: Synergizing reasoning and acting in language models. In *The eleventh international*  
625 *conference on learning representations*, 2022b. 5, 6, 21, 22
- 626 yt-dlp. yt-dlp: A feature-rich command-line audio/video downloader. [https://github.com/](https://github.com/yt-dlp/yt-dlp)  
627 [yt-dlp/yt-dlp](https://github.com/yt-dlp/yt-dlp), 2025. Accessed: 2025-12-19. 19  
628
- 629 Siyu Yuan, Kaitao Song, Jiangjie Chen, Xu Tan, Dongsheng Li, and Deqing Yang. Evoagent:  
630 Towards automatic multi-agent generation via evolutionary algorithms. In *Proceedings of the*  
631 *2025 Conference of the Nations of the Americas Chapter of the Association for Computational*  
632 *Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 6192–6217, 2025. 1, 2,  
633 7, 22
- 634 Mert Yuksekgonul, Federico Bianchi, Joseph Boen, Sheng Liu, Pan Lu, Zhi Huang, Carlos Guestrin,  
635 and James Zou. Optimizing generative ai by backpropagating language model feedback. *Nature*,  
636 639(8055):609–616, 2025. 1
- 637 Jiayi Zhang, Jinyu Xiang, Zhaoyang Yu, Fengwei Teng, Xiong-Hui Chen, Jiaqi Chen, Mingchen  
638 Zhuge, Xin Cheng, Sirui Hong, Jinlin Wang, et al. Aflow: Automating agentic workflow generation.  
639 In *The Thirteenth International Conference on Learning Representations*, 2025a. 1, 2, 3, 7, 16, 22  
640
- 641 Yao Zhang, Chenyang Lin, Shijie Tang, Haokun Chen, Shijie Zhou, Yunpu Ma, and Volker Tresp.  
642 Swarmagentic: Towards fully automated agentic system generation via swarm intelligence. *arXiv*  
643 *preprint arXiv:2506.15672*, 2025b. 1, 3, 6, 7, 18, 23  
644  
645  
646  
647

648	<b>1 Introduction</b>	<b>1</b>
649		
650	<b>2 Related Work</b>	<b>2</b>
651		
652		
653	<b>3 Mem<sup>2</sup>Evolve</b>	<b>3</b>
654	3.1 Dual-Memory Mechanism . . . . .	3
655		
656	3.2 Forward Inference . . . . .	5
657		
658	3.3 Backward Evolution . . . . .	5
659		
660	<b>4 Experiments</b>	<b>6</b>
661	4.1 Experiment Setting . . . . .	6
662		
663	4.2 Main Results . . . . .	7
664		
665	<b>5 Analysis</b>	<b>8</b>
666	5.1 RQ1: Ablation Study . . . . .	8
667		
668	5.2 RQ2: Experience-Guided Asset Creation . . . . .	8
669		
670	5.3 RQ3: Single Task Self-Evolving . . . . .	9
671		
672	5.4 RQ4: Cross Tasks Self-Evolving . . . . .	9
673		
674	<b>6 Conclusion</b>	<b>9</b>
675		
676	<b>A Mem<sup>2</sup>Evolve</b>	<b>15</b>
677	A.1 Defining Continuous and Stable Evolving . . . . .	15
678		
679	A.2 Evaluation of Existing Self-Evolving Agent Frameworks . . . . .	16
680		
681	A.3 Task Planning . . . . .	18
682		
683	A.4 Tool Creation . . . . .	18
684		
685	A.5 Assets Recruitment . . . . .	19
686		
687	A.6 Execution . . . . .	21
688		
689	<b>B Experimental Details</b>	<b>22</b>
690	B.1 Baselines . . . . .	22
691		
692	B.2 Benchmarks . . . . .	23
693		
694	B.3 RQ1: Ablation Study . . . . .	24
695		
696	B.4 RQ5: Efficiency and Cost Dynamics in Self-Evolution . . . . .	24
697		
698	<b>C Prompt Template</b>	<b>26</b>
699	C.1 Task Planning . . . . .	26
700		
701	C.2 Assess Tool Need . . . . .	26
	C.3 Tool Spec Creation . . . . .	28
	C.4 Tool Creation . . . . .	29
	C.5 Agent Creation . . . . .	31

702		
703		
704		
705		
706		
707		
708		
709		
710		
711		
712		
713		
714		
715		
716		
717		
718		
719		
720		
721		
722		
723		
724		
725		
726		
727		
728		
729		
730		
731		
732		
733		
734		
735		
736		
737		
738		
739		
740		
741		
742		
743		
744		
745		
746		
747		
748		
749		
750		
751		
752		
753		
754		
755		
	C.6	React Template . . . . . 32
	C.7	LLM as a Judge . . . . . 32
	C.8	Tool Memory Generation . . . . . 34
	C.9	Success Agent Memory Generation . . . . . 35
	C.10	Failure Agent Memory Generation . . . . . 36
	<b>D</b>	<b>Case Study</b> . . . . . <b>38</b>
	D.1	Tool Implementation for <i>Simulate Piston Platform Game</i> . . . . . 38
	D.2	Tool Implementation for <i>Youtube Audio Transcriber</i> . . . . . 40
	D.3	Experience Guidance Tool Creation . . . . . 43
	D.4	Comparison of Tool Generation With and Without Experience Guidance . . . . . 48

## A MEM<sup>2</sup>EVOLVE

### A.1 DEFINING CONTINUOUS AND STABLE EVOLVING

[BACK TO TOC](#)

We define two fundamental characteristics of self-evolving agents: **(1)** the ability to continuously evolve with minimal human intervention by persistently expanding their capabilities to solve unseen tasks, and **(2)** the capability to efficiently leverage past experience, enabling correct experience transfer and effective error avoidance for seen tasks, either offline or online.

- Optimization.** An agent system should be capable of automatically optimizing its internal instructions and coordination workflows based on environmental feedback, thereby achieving optimal task-specific performance. Traditional agent development paradigms rely heavily on manual prompt engineering or predefined interaction protocols. Such approaches are not only labor-intensive but also poorly suited to dynamically evolving task requirements. Therefore, an effective self-evolution framework should emulate a form of “backpropagation” mechanism, continuously refining agent role definitions, prompting strategies, and even multi-agent collaboration topologies through feedback signals or textual gradients derived from task execution outcomes. Crucially, this optimization should not be limited to transient runtime adjustments but should fundamentally enhance the system’s intrinsic competence in handling similar tasks.
- Experience Persistence.** To enable genuine lifelong learning, the framework must be able to transform both successful strategies and failure cases from historical tasks into long-term memory assets that persist beyond the lifecycle of a single task. Many existing methods reset system states after task completion, forcing agents to re-explore from scratch when encountering similar scenarios. This not only wastes computational resources but also allows recurring errors. Hence, a cross-task experience persistence mechanism is essential. Whether implemented via explicit external databases that store reasoning trajectories or via implicit knowledge internalization through parameter or prompt updates, this mechanism should enable rapid retrieval and reuse of prior knowledge when facing new tasks, thereby mitigating the cold-start problem and avoiding known pitfalls.
- Agent Creation.** The framework should not depend on predefined expert agent modules with fixed roles, prompts, or decision logics. Instead, it should be capable of dynamically constructing optimal expert agent teams conditioned on task demands. This capability is critical for high-level, complex planning tasks, where increasing task complexity typically entails decomposing the problem into multiple subtasks with distinct objectives. A single general-purpose agent is often insufficient to handle all subtasks effectively, necessitating specialized agents that each address the components they are best suited for. However, given the vast diversity of real-world tasks, manually predefined expert agents cannot cover all possible scenarios. The system must therefore support dynamic, task-driven agent generation.
- Tool Creation.** By invoking tools, agent systems can substantially expand their capability boundaries and overcome the limitations imposed by static knowledge. Tools enable access to real-time information, execution of complex mathematical reasoning, and completion of specialized operations. However, task-specific tools typically require careful manual design. When confronted with general, previously unseen tasks, human developers are often still required to create new tools, which is inherently unscalable. To enable continual capability expansion, the framework must therefore possess the ability to autonomously generate tools.
- Experience-Guided Creation.** When encountering unseen tasks that require the generation of new agents or tools, the framework should leverage its internal assets and accumulated memory to guide the creation process. For example, when a new task involves parsing YouTube video subtitles, previously generated tool documentation for downloading YouTube videos can serve as a reference to facilitate new tool construction. This experience-guided mechanism improves the stability of generated tools and agents, reduces randomness and hallucinations in large language model outputs, and thereby enables a more robust and reliable evolutionary process.

## A.2 EVALUATION OF EXISTING SELF-EVOLVING AGENT FRAMEWORKS

[BACK TO TOC](#)

- 810

811

812 • **DyLAN** (Liu et al., 2023) models multi-agent collaboration as a Temporal Feed-Forward

813 Network, implementing a "Team Optimization" stage that utilizes a backward message-

814 passing mechanism to calculate "Agent Importance Scores" based on unsupervised peer

815 ratings. **DyLAN satisfies Optimization:** it actively employs environmental feedback

816 to refine the collaboration topology iteratively. By identifying and selecting the most

817 contributory agents while pruning low-performing ones, DyLAN automatically optimizes

818 the team composition and interaction structure for specific tasks, aligning with the definition

819 of optimizing collaboration logic and topology. However, **DyLAN fails Agent Creation**

820 **and Tool Creation:** it does not generate new expert definitions or tools from scratch;

821 instead, it relies on selecting the best subset from a fixed, pre-defined candidate pool of

822 agents. Finally, it **offers limited Experience Persistence:** while it can reuse calculated

823 importance scores for similar tasks, it lacks a semantic memory mechanism to guide the

824 generation of new assets for entirely unseen domains.
- 825 • **DSPy** (Khatab et al., 2023) introduces a programming model that abstracts language model

826 pipelines as text transformation graphs, allowing users to define declarative modules (e.g.,

827 ChainOfThought) via natural language signatures. **DSPy satisfies Optimization:** it employs

828 a compiler with various "teleprompters" to automatically refine the pipeline's instructions

829 or fine-tune the underlying language model weights based on metric-driven feedback and

830 bootstrapped demonstrations. However, **DSPy fails Agent Creation and Tool Creation:**

831 it relies on the user to explicitly define the program structure, the flow of control, and the

832 specific modules and tools to be used, rather than autonomously synthesizing new agent roles

833 or executable tools from scratch. Furthermore, **DSPy fails Experience Persistence:** it treats

834 experience utilization as a discrete compilation process rather than a continuous memory

835 accumulation. Once the pipeline is compiled, the historical traces are frozen into static

836 few-shot examples or weights, lacking a dynamic, retrievable memory bank to persistently

837 store and reuse new inference trajectories for future tasks.
- 838 • **ReasoningBank** (Ouyang et al., 2025) introduces a memory framework that distills general-

839 izable reasoning strategies from both successful and failed trajectories, enabling agents to

840 retrieve relevant insights for new tasks. **ReasoningBank satisfies Experience Persistence:**

841 it explicitly stores abstracted reasoning patterns in a long-term memory bank, allowing the

842 system to mitigate the cold-start problem by transferring knowledge across tasks and pre-

843 venting the repetition of past errors. However, **ReasoningBank fails Optimization:** while

844 it improves performance via RAG, it does not structurally optimize the agent's topology,

845 internal prompt templates, or parameters. The agent's core functionality remains static,

846 relying on external memory injection rather than internal refinement. It likewise fails **Agent**

847 **Creation and Tool Creation**, as it operates with a fixed agent architecture (e.g., ReAct),

848 leveraging dynamic memory, rather than generating new agent entities or executable tools

849 from scratch. Consequently, it also fails **Experience-Guided Creation**, as there is no asset

850 generation process to be guided by its rich memory.
- 851 • **AFlow** (Zhang et al., 2025a) reformulates agentic system development as a search problem

852 over code-represented workflows, utilizing Monte Carlo Tree Search (MCTS) to iteratively

853 explore and optimize the design space. **AFlow satisfies Optimization:** it treats the workflow

854 structure and node-level prompts as hyperparameters, optimizing them based on execution

855 feedback (e.g., success rates, costs) to find the most effective graph topology. It also

856 satisfies **Agent Creation and Experience-Guided Creation:** the framework autonomously

857 constructs new workflow nodes (effectively new agents) and connections by leveraging the

858 search history (MCTS values) to guide the generation process, moving away from manual

859 engineering. However, **AFlow fails Tool Creation:** it focuses on orchestrating the flow of

860 LLM calls and existing tools rather than synthesizing new executable tool code from scratch.

861 Furthermore, it fails **Experience Persistence:** the experience is utilized only during the

862 offline search phase to produce a static, compiled workflow; it lacks a dynamic, long-term

863 memory mechanism to continuously accumulate and retrieve reasoning patterns for lifelong

learning across different tasks.
- 864 • **AgentSquare** (Shang et al., 2025) proposes a search-based framework that automates the

design of LLM agents by exploring a modular space comprising planning, reasoning, tool

use, and memory modules. **AgentSquare satisfies Optimization:** it treats the agent design

process as an objective function maximization problem, iteratively refining the agent’s architecture (i.e., module combinations) based on evaluation feedback to find the optimal configuration. It also satisfies **Agent Creation** and **Experience-Guided Creation**: the framework autonomously synthesizes new agent instances by recombining modules and utilizes an “Experience Pool” (containing history of evaluated agent-task pairs) to train a surrogate model, which efficiently guides the generation of new candidates towards high-performance regions. It further satisfies **Experience Persistence** by explicitly storing these search trajectories and evaluation results in the Experience Pool, enabling the system to learn from past search iterations. However, **AgentSquare fails Tool Creation**: while it optimizes the *mechanism* of tool usage (e.g., choosing between ReAct or Plan-and-Solve), it relies on orchestrating existing tools rather than generating new executable tool implementations from scratch to extend capabilities.

- **ANN** (Ma et al., 2025) conceptualizes multi-agent systems as neural networks, treating agents as learnable nodes and their communication as edges. **ANN satisfies Optimization**: drawing inspiration from backpropagation, it introduces a “textual backpropagation” mechanism that computes textual gradients based on error feedback to iteratively update the agents’ system prompts (which function as learnable weights). However, **ANN fails Agent Creation and Tool Creation**: the framework operates on pre-defined architectural topologies (e.g., Chain, Stack, or Grid structures) with a fixed number of agent nodes; it optimizes the *behavior* of these existing agents rather than autonomously synthesizing new agent roles or executable tools from scratch. Furthermore, it fails **Experience Persistence**: similar to traditional model training, the learned experience is implicitly internalized into the optimized prompt parameters during the optimization phase. It lacks a dynamic, explicitly retrievable memory bank to persistently store reasoning trajectories, limiting its ability to support lifelong learning or transfer insights to entirely new domains without re-training. Consequently, it also fails **Experience-Guided Creation**.
- **Alita** (Qiu et al., 2025) introduces a self-evolving framework that enables an LLM agent to dynamically expand its capability boundaries by creating and integrating new tools. **Alita satisfies Tool Creation**: it employs a “Creator” module that autonomously synthesizes executable Python tools from scratch to address tasks where existing tools are insufficient. It also satisfies **Optimization** and **Experience Persistence**: the framework maintains an explicit “Experience Pool” of successful tool-use trajectories and utilizes a “Promptist” module to retrieve relevant demonstrations and refine the agent’s prompts based on execution feedback. However, **Alita fails Agent Creation**: it operates as a single-agent system that evolves its tool library, rather than synthesizing new agent roles or collaborative teams. Furthermore, it fails **Experience-Guided Creation** (in the context of asset generation): while it uses experience to optimize *usage* prompts, the generation of *new tools* is driven by immediate task failures and reflection, without leveraging a retrieval mechanism over historical creation artifacts to guide the synthesis process.
- **ToolMaker** (Wölflein et al., 2025) introduces a dual-LLM framework where a “Tool Maker” autonomously generates reusable Python tools (functions) to address specific tasks, which are then utilized by a “Tool User” for subsequent problem-solving. **ToolMaker satisfies Tool Creation**: its core mechanism is the synthesis of executable code to encapsulate reasoning steps into reusable tools, thereby explicitly expanding the agent’s action space. It also satisfies **Optimization**: during the tool creation phase, it employs a verification loop (utilizing unit tests) to validate the generated code and uses error feedback to iteratively refine and debug the tool until it functions correctly. However, **ToolMaker fails Agent Creation**: the framework operates with fixed, pre-defined roles (Maker and User) rather than synthesizing new agent personas or collaborative team structures from scratch. Furthermore, it fails **Experience Persistence** (in the context of cross-task learning) and **Experience-Guided Creation**: while the generated tools are stored for reuse on instances of the *same* task, the framework does not maintain a retrievable memory bank of creation strategies or past artifacts to guide the generation process for entirely *new, unseen* tasks, effectively facing the cold-start problem for each new domain.
- **AgentVerse** (Chen et al., 2023) proposes a flexible multi-agent framework that orchestrates the problem-solving process through four key stages: Expert Recruitment, Collaborative Decision Making, Action Execution, and Evaluation. **AgentVerse satisfies Agent Creation**: utilizing the Expert Recruitment mechanism, the framework autonomously generates and

918 customizes new agent roles and descriptions tailored to the specific progress of the task,  
 919 rather than relying on a fixed set of pre-defined personas. It also satisfies **Optimization**: the  
 920 Evaluation stage provides real-time feedback on the agents’ outcomes, which is used to iter-  
 921 atively refine the collaborative decision-making process and adjust the team’s composition  
 922 or strategies during runtime. However, **AgentVerse fails Tool Creation**: while agents can  
 923 utilize existing tools or execute code, the framework focuses on evolving the *team struc-*  
 924 *ture* and *agent personas* rather than synthesizing new, reusable executable tool definitions  
 925 from scratch to expand the action space. Furthermore, it fails **Experience Persistence** and  
 926 **Experience-Guided Creation**: the optimization is confined to the immediate context of the  
 927 current task loop; it lacks a long-term, retrievable memory mechanism to store successful  
 928 collaboration patterns or reasoning trajectories for future cross-task transfer, meaning each  
 929 new task session effectively starts without historical guidance.

- 930 • **AutoAgents** (Chen et al., 2024) introduces an automatic agent generation framework that  
 931 dynamically synthesizes a team of specialized agents (including roles and expert profiles)  
 932 tailored to the specific input task. **AutoAgents satisfies Agent Creation**: it leverages a  
 933 “Planner” to decompose the task and autonomously generate the identities and descriptions  
 934 of the necessary expert agents, rather than selecting from a pre-defined library. It also  
 935 satisfies **Optimization**: it employs an “Observer” mechanism (Agent Observer and Plan  
 936 Observer) to review the generated agents and execution plans, providing feedback to refine  
 937 and optimize the team structure and workflows before execution. However, **AutoAgents**  
 938 **fails Tool Creation**: while the generated agents can utilize existing tools, the framework  
 939 focuses on synthesizing the *agents* themselves, not generating new executable tool code from  
 940 scratch to expand the system’s capabilities. Furthermore, it fails **Experience Persistence** and  
 941 **Experience-Guided Creation**: the generation process is effectively zero-shot for each new  
 942 task instance; it does not maintain a long-term, retrievable memory of past successful agent  
 943 configurations or planning trajectories to guide the generation of future agents, tackling each  
 944 task as an isolated event without accumulation of experience.
- 945 • **SwarmAgentic** (Zhang et al., 2025b) applies Swarm Intelligence (SI) principles (specifically  
 946 Particle Swarm Optimization) to the domain of agent system design, treating agents and  
 947 tools as particles that evolve in a search space. **SwarmAgentic satisfies Agent Creation**  
 948 **and Tool Creation**: the framework autonomously synthesizes both the agent definitions  
 949 (roles, prompts) and executable tool code (Python functions) from scratch to construct a  
 950 functional system, rather than selecting from a fixed library. It also satisfies **Optimization**:  
 951 it utilizes a velocity-based update mechanism to iteratively refine the agents’ prompts and  
 952 tools based on the “personal best” and “global best” feedback found during the swarm search  
 953 process. However, **SwarmAgentic fails Experience Persistence**: the optimization history and  
 954 learned patterns are transient, utilized only to converge on a solution for the current task  
 955 instance. It does not maintain a persistent, retrievable memory bank of successful designs  
 956 to support lifelong learning across different tasks. Consequently, it also fails **Experience-**  
 957 **Guided Creation**, as the initialization of new systems for unseen tasks relies on zero-shot  
 958 generation rather than being informed by a repository of historical assets.

### 958 A.3 TASK PLANNING

[BACK TO TOC](#)

959 In real-world settings, tasks are typically accomplished through multiple steps. To improve the  
 960 specialization of subsequently created tools and to fully leverage expert agents by assigning them  
 961 distinct roles aligned with their respective strengths, the system first decomposes a complex task into  
 962 a set of sub-tasks during the planning stage. Each sub-task specifies its objective, the expected output  
 963 format, and its dependencies on other sub-tasks—namely, the results from prerequisite sub-tasks  
 964 required for execution.

### 966 A.4 TOOL CREATION

[BACK TO TOC](#)

968 When a sub-task exceeds the agent’s existing capability scope, the system initiates a tool-creation  
 969 workflow to expand its capability boundaries. However, synthesizing effective tools based solely on  
 970 sub-task descriptions proves inadequate. Our empirical analysis identifies three key limitations: (1)  
 971 brief sub-task descriptions provide insufficient constraints, resulting in unstable and inconsistent tool  
 generation; (2) tools derived solely from the model’s internal, static knowledge often lack practical

972 usability; and (3) the inherent stochasticity of model outputs hinders the reproducibility of successful  
 973 tool-generation processes and prevents effective reuse of failure cases.

974 To overcome these challenges, we introduce a three-stage tool synthesis strategy: (1) tool specification  
 975 generation to formalize functionality and interfaces; (2) tool documentation and experience collection  
 976 to ground tool usage and accumulate actionable knowledge; and (3) tool implementation to produce  
 977 reliable and reusable tools.  
 978

- 979 • **Tool Spec Generation.** Inspired by specification-driven development paradigms (GitHub,  
 980 2025), we require the model to first generate a formal tool specification before implementing  
 981 the tool itself. As illustrated in Figure 4, this specification includes the *tool name*, a *concise*  
 982 *description*, *input parameters*, *output format*, and *core logic*. The core logic is articulated  
 983 as a sequence of concrete steps that explicitly describe the tool’s execution process (e.g.,  
 984 “Step 1: validate input compliance”). By introducing this intermediate specification stage,  
 985 the agent can generate tools by directly adhering to well-defined requirements, thereby  
 986 improving controllability, consistency, and overall accuracy in the tool creation process.
- 987 • **Tool Documentation and Experience Collection.** To further mitigate the limitation of  
 988 relying solely on the parametric, static knowledge, we incorporate an additional grounding  
 989 step prior to tool generation. Specifically, the agent leverages a web search tool to retrieve  
 990 external documentation, such as open-source tool descriptions on GitHub<sup>7</sup> and debugging  
 991 discussions from Stack Overflow<sup>8</sup>. In parallel, the agent queries its Experience Memory  
 992 using the generated tool specification to retrieve relevant development references. For  
 993 example, both a tool for extracting basic YouTube video metadata and a tool for downloading  
 994 YouTube videos can be grounded in documentation for the open-source utility *yt-dlp* (yt-dlp,  
 995 2025). By integrating externally sourced documents with experience-based retrieval, the  
 996 tool creation process is better grounded in real-world implementations, leading to more  
 997 practical and reliable tools.
- 998 • **Tool Implementation.** With a well-defined tool specification and sufficiently rich tool  
 999 documentation and implementation experience stored in memory, the system can proceed to  
 1000 generate the corresponding tool. As shown in Code 1, we intentionally encapsulate each tool  
 1001 in an MCP-compliant format, ensuring that it can be seamlessly integrated by different LLM  
 1002 backbones. This design enables model-agnostic interoperability and facilitates efficient  
 1003 reuse in subsequent applications.

## 1004 A.5 ASSETS RECRUITMENT

[BACK TO TOC](#)

1005 To optimize the utilization of tools and expert agents stored in Asset Memory, Mem<sup>2</sup>EvoImple-  
 1006 ments an *Assets Recruitment* phase prior to task execution. This mechanism operates at the granularity  
 1007 of sub-tasks. Let  $\mathbf{q} = \text{embedding}(s_i)$  denote the embedding of the current sub-task description.

1008 **Expert Agent Retrieval** We query the agent memory  $\mathcal{M}_{agt}$  to find the most proficient expert. The  
 1009 retrieval key for a candidate agent  $a_i$  is defined as  $\mathbf{h}_{a_i} = \text{embedding}(\rho_i \oplus \epsilon_i \oplus \sigma_i)$ , derived from  
 1010 its role, expertise, and behavior suggestions. We select the optimal agent  $a^*$  by retrieving the Top-1  
 1011 candidate that exceeds a similarity threshold  $\delta$ :  
 1012

$$1013 a^* = \arg \max_{a_i \in \mathcal{M}_{agt}} \left\{ \cos(\mathbf{q}, \mathbf{h}_{a_i}) \mid \cos(\mathbf{q}, \mathbf{h}_{a_i}) > \delta \right\} \quad (12)$$

1014 If the set is empty, a new agent generation process is triggered.  
 1015

1016 **Tool Retrieval** Similarly, for tool memory  $\mathcal{M}_{tool}$ , the key is  $\mathbf{h}_{t_j} = \text{embedding}(n_j \oplus d_{func,j})$ . To  
 1017 balance functional support with context window constraints, we retrieve the Top- $K$  relevant tools to  
 1018 form the available toolset  $\mathbb{T}_{avail}$ :  
 1019

$$1020 \mathbb{T}_{avail} = \text{Top-K}_{t_j \in \mathcal{M}_{tool}} \left\{ t_j \mid \cos(\mathbf{q}, \mathbf{h}_{t_j}) > \delta \right\} \quad (13)$$

1024 <sup>7</sup><https://github.com>

1025 <sup>8</sup><https://stackoverflow.com/>

```

1026 [
1027   {
1028     "tool_name": "simulate_piston_platform_game",
1029     "tool_description": "Simulates a specific ping-pong game mechanism involving a
1030     ball queue, a limited-capacity platform, and three random pistons with complex
1031     ejection/replacement rules. Used to calculate the probability of each ball number
1032     being 'ejected by a piston' (winning) versus being 'released' (eliminated).",
1033     "input_parameters": [
1034       {
1035         "name": "num_balls",
1036         "type": "integer",
1037         "description": "Total number of balls in the queue, numbered 1 to N.",
1038         "default": 100
1039       },
1040       {
1041         "name": "num_simulations",
1042         "type": "integer",
1043         "description": "Monte Carlo iterations.",
1044         "default": 100000
1045       }
1046     ],
1047     "output_format": {
1048       "type": "object",
1049       "properties": {
1050         "win_probabilities": {
1051           "type": "object",
1052           "description": "Mapping of ball number to its probability of being
1053           ejected by a piston (Winning).",
1054         },
1055         "best_choice": {
1056           "type": "integer",
1057           "description": "The ball number with the highest win probability."
1058         }
1059       }
1060     },
1061     "core_logic": [
1062       "Step 1: Initialize `win_counts` for all balls to 0.",
1063       "Step 2: Run loop `num_simulations` times.",
1064       "Step 3: In each sim, initialize a queue `deck` [1..num_balls] and a `
1065       platform` holding the first 3 balls.",
1066       "Step 4: While platform is not empty, randomly select a piston (1, 2, or 3)
1067       with equal probability.",
1068       "Step 5: Apply Piston Rules:",
1069       "  - If Piston 1: Eject Pos 1 (Win). Move Pos 2->1, Pos 3->2. Refill 1 ball
1070       from deck to Pos 3.",
1071       "  - If Piston 2: Eject Pos 2 (Win). Release Pos 1 (Loss/Die). Move Pos 3->
1072       1. Refill 2 balls from deck to Pos 2 & 3.",
1073       "  - If Piston 3: Eject Pos 3 (Win). Release Pos 1 (Loss/Die). Move Pos 2->
1074       1. Refill 2 balls from deck to Pos 2 & 3.",
1075       "Step 6: If a ball is 'Ejected' (Win), increment its count in `win_counts`.
1076       If 'Released', do nothing.",
1077       "Step 7: Continue until platform and deck are empty.",
1078       "Step 8: Calculate probabilities = wins / total_simulations."
1079     ]
1080   }
1081 ]

```

Figure 4: **Specification of the tool for simulating the Piston Platform game.** The specification includes the tool name and description, detailed definitions of input parameters and output formats—where each parameter is characterized by its name, type, description, and default value—as well as the core logic of the tool implementation, guiding subsequent tool creation.

```

1080
1081 {
1082     "role": "Probability Simulation Analyst",
1083     "expertise": "Specializes in stochastic modeling and quantitative analysis to
1084     derive probabilities from complex mechanical simulations.",
1085     "suggestions": [
1086         "Execute multiple simulation runs to ensure statistical significance of
1087         the results.",
1088         "Aggregate ejection data to calculate the specific probability for each
1089         ball.",
1090         "Identify the ball with the maximum ejection frequency from the dataset
1091         ."
1092     ],
1093     "tools": [
1094         "simulate_ping_pong_game"
1095     ]
1096 }
1097

```

Figure 5: **Specification of the probabilistic simulation expert.** The specification defines the expert agent’s role, areas of expertise, suggested strategies or recommendations, and the list of tools available for use during task execution.

## A.6 EXECUTION

[BACK TO TOC](#)

Since expert agents must frequently invoke tools to interact with the external environment and make subsequent decisions based on environmental feedback, we adopt the ReAct (Yao et al., 2022b) framework, which alternates among **Think**, **Action**, and **Observation** steps. Specifically, we design a standardized ReAct system prompt template with dedicated placeholders for prerequisite results from dependent sub-tasks, the expert role, task-specific suggestions, and the set of available tools. During execution, these placeholders are dynamically populated with agent-specific content for each expert agent. This templated design ensures both the stability of the reasoning–action loop and the extensibility of the framework across different expert roles and task settings.

## B EXPERIMENTAL DETAILS

### B.1 BASELINES

[BACK TO TOC](#)

In this section, we provide detailed descriptions of the baseline frameworks employed in our evaluation, categorized by their operational paradigms.

#### Naive Large Language Models

- **Direct:** This is the most fundamental approach, where the task description is fed directly into the Large Language Model (LLM) without any intermediate reasoning steps or external tools. It serves as a baseline to measure the inherent zero-shot capability of the model.
- **CoT (Wei et al., 2022):** CoT enhances the reasoning capabilities of LLMs by prompting them to generate a series of intermediate reasoning steps before producing the final answer. This method is particularly effective for complex tasks requiring multi-step logic.
- **ReAct (Yao et al., 2022b):** ReAct synergizes reasoning and acting by allowing the model to generate reasoning traces and task-specific actions (such as web searches) in an interleaved manner. This enables the agent to interact with external environments to retrieve information and update its context dynamically.
- **OpenAI Deep Research:** A commercial-grade autonomous research agent developed by OpenAI. It is designed to perform deep, multi-step research tasks by browsing the web, synthesizing information from multiple sources, and generating comprehensive reports, representing the state-of-the-art in proprietary agent systems.

#### Experience-Centric Frameworks

- **DyLAN (Liu et al., 2023):** DyLAN is a framework that models multi-agent collaboration as a Temporal Feed-Forward Network. It introduces a “Team Optimization” stage utilizing a backward message-passing mechanism to calculate Agent Importance Scores. By actively identifying high-contribution agents and pruning low-performing ones, DyLAN iteratively optimizes the team’s collaboration topology, though it relies on a fixed pool of agents rather than creating new ones.
- **EvoAgent (Yuan et al., 2025):** EvoAgent applies evolutionary algorithms to multi-agent systems, treating agents as individuals in a population. It employs operations such as crossover and mutation on agent prompts to iteratively evolve their behaviors. This allows the system to discover more effective agent personas and strategies over time without manual prompt engineering.
- **AFlow (Zhang et al., 2025a):** AFlow reformulates agentic system development as a search problem over code-represented workflows. Utilizing Monte Carlo Tree Search (MCTS), it iteratively explores and optimizes the design space of workflow structures and node-level prompts. AFlow autonomously constructs new workflow nodes and connections guided by search history, moving away from manual engineering to find the most effective graph topology for specific tasks.
- **DSPy (Khatab et al., 2023):** DSPy introduces a programming model that abstracts LM pipelines as text transformation graphs. It employs a compiler with various “teleprompters” to automatically refine the pipeline’s instructions or fine-tune the underlying language model weights based on metric-driven feedback. While it optimizes the pipeline effectively, it relies on user-defined program structures rather than autonomously synthesizing new agent roles or tools.

#### Capacity-Centric Frameworks (Tool-Generative)

- **Alita<sup>9</sup> (Qiu et al., 2025):** Alita is a self-evolving framework designed to dynamically expand an agent’s capability boundaries. It features a “Creator” module that autonomously

<sup>9</sup>We use the implementation at <https://github.com/ryantzr1/OpenAlita> as the official code is unavailable.

synthesizes executable Python tools from scratch to address tasks where existing tools are insufficient. Additionally, Alita utilizes a “Promptist” module and an explicit Experience Pool to refine usage prompts based on execution feedback, enabling continuous adaptation.

- **ToolMaker** (Wölflein et al., 2025): ToolMaker adopts a dual-LLM framework comprising a “Tool Maker” and a “Tool User.” The Maker autonomously generates reusable Python tools (functions) to encapsulate reasoning steps for specific tasks, while the User applies them for problem-solving. It includes a verification loop with unit tests to ensure the reliability of the generated code, explicitly expanding the agent’s action space through code synthesis.

### Capacity-Centric Frameworks (Agent-Generative)

- **AgentVerse** (Chen et al., 2023): AgentVerse proposes a flexible multi-agent framework orchestrating the problem-solving process through Expert Recruitment, Collaborative Decision Making, Action Execution, and Evaluation. It autonomously generates and customizes new agent roles tailored to the task progress. The framework uses real-time feedback to refine the collaborative process and adjust team composition during runtime.
- **AutoAgents** (Chen et al., 2024): AutoAgents introduces an automatic agent generation framework that dynamically synthesizes a team of specialized agents tailored to the input task. Leveraging a “Planner” to decompose tasks, it autonomously generates expert identities and descriptions. An “Observer” mechanism further reviews and refines the execution plans and team structure, ensuring the generated agents are optimized for the specific problem context.
- **SwarmAgentic** (Zhang et al., 2025b): SwarmAgentic applies Swarm Intelligence principles, specifically Particle Swarm Optimization (PSO), to agent system design. It treats agents and tools as particles evolving in a search space, autonomously synthesizing both agent definitions and executable tool code. The framework utilizes velocity-based updates to iteratively refine prompts and tools based on “personal best” and “global best” feedback found during the swarm search.

## B.2 BENCHMARKS

[BACK TO TOC](#)

Table 5: **Overview of the benchmarks, domains, test set sizes, and evaluation metrics used in experiments.** † indicates that the test set was randomly sampled. ‡Average Score is the mean of Delivery Rate, Micro/Macro Commonsense Constraint Pass Rate, and Micro/Macro Hard Constraint Pass Rate.

Benchmark	Domain	Test Size	Metric
GAIA (Mialon et al., 2024)	General Assistant	166	Pass@1
ALFWorld (Shridhar et al., 2020)	Embodied Task	134	Success Rate
HotpotQA (Yang et al., 2018)	Multi-hop QA	500 <sup>†</sup>	Exact Match (EM)
2WikiMultiHopQA (Ho et al., 2020)	Multi-hop QA	500 <sup>†</sup>	Exact Match (EM)
AIME 24	Math Reasoning	30	Pass@1
AIME 25	Math Reasoning	30	Pass@1
TravelPlanner (Xie et al., 2024)	Complex Planning	1,000	Average Score <sup>‡</sup>
WebShop (Yao et al., 2022a)	Web Navigation	251	Success Rate

To evaluate the general-purpose task-solving capability of DEMA, we conduct experiments across six task categories and eight benchmarks, as summarized in Table 5.

- **GAIA**: GAIA is a benchmark designed to assess the capabilities of general-purpose AI assistants. It consists of 466 real-world, scenario-based questions covering daily tasks, scientific reasoning, web browsing, and tool usage. While these tasks are conceptually simple for humans, they remain challenging for advanced AI systems. We report Pass@1 as the primary evaluation metric.
- **ALFWorld**: ALFWorld aligns text-based games with embodied environments to evaluate an agent’s ability to reason and act in interactive settings. It requires the agent to understand

high-level goals and execute a sequence of low-level actions to interact with objects. We evaluate our method on 134 tasks and use Success Rate as the evaluation metric.

- **HotpotQA**: HotpotQA is a question-answering benchmark that challenges agents to perform multi-hop reasoning across multiple documents to locate relevant facts. In our experiments, we randomly sample a subset of 500 instances as the test set. Performance is evaluated using the Exact Match (EM) metric.
- **2WikiMultihopQA**: Similar to HotpotQA, 2WikiMultihopQA evaluates multi-hop reasoning capabilities over Wikipedia articles and features complex queries that require synthesizing information from multiple sources. We randomly sample 500 instances for testing and use EM for evaluation.
- **AIME 24/25**: These benchmarks correspond to the problem sets from the 2024 and 2025 American Invitational Mathematics Examinations. Each dataset consists of 30 high-difficulty problems that test advanced mathematical reasoning and creative problem-solving abilities. We use Pass@1 to measure accuracy.
- **TravelPlanner**: TravelPlanner evaluates language agents in complex tool-use and long-horizon planning scenarios, such as generating travel itineraries under strict constraints. We use the full test set of 1,000 instances. The final score is computed as the average of five metrics: Delivery Rate, Micro Commonsense Constraint Pass Rate, Macro Commonsense Constraint Pass Rate, Micro Hard Constraint Pass Rate, and Macro Hard Constraint Pass Rate.
- **WebShop**: WebShop is a simulated e-commerce environment containing over one million real-world products. It tests an agent’s ability to navigate web pages, search, browse, and select options to fulfill user instructions. We evaluate on 251 test instances and measure performance using Success Rate.

B.3 RQ1: ABLATION STUDY

[BACK TO TOC](#)

Table 6: **Ablation study of Mem<sup>2</sup>Evolve**. Pass@1 scores are reported. Performance drops (↓) relative to the full model are shown in parentheses.

Method	GAIA	Embodied	Multi-Hop QA		Math		Planning	Web	Avg.
	Total	ALFWorld	HotpotQA	2Wiki	AIME24	AIME25	TravelPlanner	WebShop	
DMEA (Ours)	76.31	94.31	60.80	82.00	76.70	73.33	59.25	39.20	70.24
<i>w/o Asset Creation</i>									
w/o Tool Creation	21.69 (↓ 54.62)	94.10 (↓ 0.21)	59.50 (↓ 1.30)	81.50 (↓ 0.50)	66.67 (↓ 10.03)	60.00 (↓ 13.33)	57.15 (↓ 2.10)	39.05 (↓ 0.15)	59.96 (↓ 10.28)
w/o Expert Agent Creation	75.30 (↓ 1.01)	93.65 (↓ 0.66)	59.00 (↓ 1.80)	81.00 (↓ 1.00)	73.33 (↓ 3.37)	70.00 (↓ 3.33)	57.08 (↓ 2.17)	38.80 (↓ 0.40)	68.52 (↓ 1.72)
<i>w/o Experience Distillation</i>									
w/o Tool Memory	67.47 (↓ 8.84)	94.25 (↓ 0.06)	60.00 (↓ 0.80)	81.50 (↓ 0.50)	70.00 (↓ 6.70)	66.67 (↓ 6.66)	57.85 (↓ 1.40)	39.15 (↓ 0.05)	67.11 (↓ 3.13)
w/o Agent Memory	74.70 (↓ 1.61)	88.06 (↓ 6.25)	56.80 (↓ 4.00)	77.40 (↓ 4.60)	73.33 (↓ 3.37)	70.00 (↓ 3.33)	49.36 (↓ 9.89)	34.40 (↓ 4.80)	65.51 (↓ 4.73)

B.4 RQ5: EFFICIENCY AND COST DYNAMICS IN SELF-EVOLUTION

[BACK TO TOC](#)

In this section, we investigate the dynamics of task completion efficiency during the self-evolution process. Specifically, we analyze the variation in token consumption across the entire forward inference phase, encompassing planning, creation, debugging, and execution, as the number of executed tasks increases under a single-task setting without memory initialization. To rigorously assess this, we organize the GAIA tasks by increasing difficulty (sequentially from Level 1 to Level 3) and compare Mem<sup>2</sup>Evolve against a w/o Dual-Memory baseline, where the agent is forced to reset to 0% memory before every task, thereby precluding the persistent storage of generated assets and distilled experiences.

As illustrated in Figure 6, while the baseline exhibits step-like cost increases strictly proportional to task difficulty with no visible optimization, Mem<sup>2</sup>Evolve demonstrates a clear evolutionary efficiency advantage. First, within each difficulty level, the system shows a continuous downward trend in token consumption, confirming that as Asset Memory expands, the agent increasingly shifts

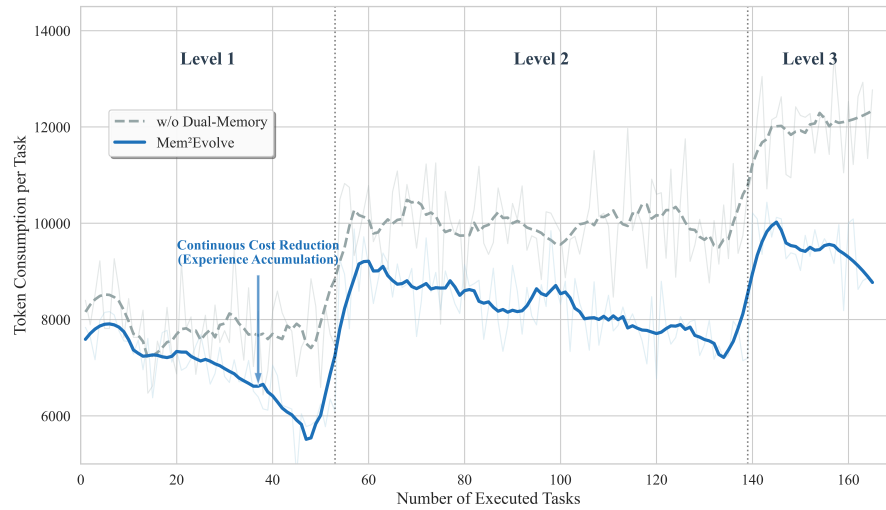


Figure 6: **Evolutionary cost dynamics on the GAIA.** As the number of executed tasks increases, Mem<sup>2</sup>Evolve demonstrates a continuous downward trend in token consumption compared to the w/o Dual-Memory. This indicates that by retrieving and reusing accumulated assets and experience from the Dual-Memory, the agent significantly reduces the overhead of exploration and creation, leading to progressively higher efficiency even as task difficulty increases.

from an expensive creation mode to a cost-effective reuse mode. Second, when transitioning to harder capability levels (e.g., from Level 1 to Level 2), Mem<sup>2</sup>Evolve utilizes prior generalizable experiences to significantly dampen the cost spike observed in the baseline; this indicates that the system successfully transfers knowledge from simpler tasks to lower the cold-start cognitive load for more complex problems, establishing that Mem<sup>2</sup>Evolve becomes progressively faster and more efficient as it evolves.

## C PROMPT TEMPLATE

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

### Prompt for Task Planning

[Back to ToC](#)

You are a task planning expert. Analyze the query carefully and break down complex queries into logical, actionable steps.

USER QUERY: {query}

You MUST output EXACTLY ONE JSON object with this structure:

```
{
  "sub_tasks": [
    {
      "description": "Specific actionable task...",
      "dependencies": []
    },
    {
      "description": "Specific actionable task...",
      "dependencies": []
    }
  ]
}
```

RULES:

1. Each task description should be COMPLETE and SPECIFIC
2. Use task numbers (1,2,3...) for dependencies
3. First task usually has no dependencies: "dependencies": []
4. If task 2 depends on task 1: "dependencies": [1]
5. If task 3 depends on tasks 1 and 2: "dependencies": [1, 2]
6. Preserve the EXACT units and format requirements from the original query in task descriptions
7. Output pure JSON format with no other content.

### Prompt for Assess Tool Need

[Back to ToC](#)

You are a Principal Engineer. Your goal is to analyze the given multi-step task plan and determine the most efficient way to solve it. Your primary objective is to use LLM-Native capabilities and reuse existing tools. You will only determine that new tools are needed (need\_creation: true) as a last resort.

- Core Principles for Analysis

- Maximize LLM-Native Capabilities (Highest Priority): You (the LLM) are the primary tool. You can perform many tasks natively without code. Do NOT propose a tool for any step that involves:

- Reasoning, planning, or making decisions based on provided context.

- Extracting, reformatting, filtering, sorting, or summarizing data from context.

- Simple data transformations or list operations (e.g., finding an item, counting).

- Any task that doesn't strictly require one of the "Why Code?" principles below.

- Reuse First, Create Last (Second Priority): After exhausting LLM-Native capabilities, your next goal is to reuse available\_tools. Do NOT propose duplicates or near-duplicates.

- Analyze the Full Plan: Review the entire sequence of steps to understand overall goals and data flow.

- Deconstruct Each Step: Mentally decompose each step into atomic required\_capabilities.

```

1404
1405   - Justify the Need for a Tool (Why Code?): A step requires a code tool only if
1406   it involves at least one of:
1407     - Complex Calculations: Non-trivial math, statistics, physics (beyond simple
1408     arithmetic).
1409     - Stateful Simulation/Iteration: Executing a loop many times while tracking
1410     changing state.
1411     - Complex Data Manipulation: Creating/modifying nested data structures
1412     according to precise rules (not simple extraction).
1413     - External I/O: Accessing APIs, files, databases, or system resources.
1414     - Deterministic Logic: When a precise, non-negotiable output is required based
1415     on complex rules.
1416
1417   - Decision Procedure
1418     - Decompose Plan: Break down the plan into a flat list of atomic
1419     required_capabilities.
1420     - Evaluate Each Capability (3-step check):
1421       a. LLM-Native? Can this capability be handled directly by the LLM (based on
1422       Maximize LLM-Native Capabilities)? If yes - it's covered.
1423       b. Existing Tool? If not LLM-native - can it be handled by one or more
1424       available_tools? If yes - add those tool(s) to matching_tools.
1425       c. Gap Identified? Only if both above answers are NO - this capability is
1426       truly missing/gap.
1427
1428   - Determine Final Output
1429     If all capabilities are covered by either above method:
1430     need_creation=false; missing_tools=[]
1431     Else if any essential gaps remain from above analysis:
1432     need_creation=true; propose minimal set of missing_tools covering ONLY
1433     those gaps
1434
1435   STRICT TYPE CONTRACT
1436   - For every field listed below, emit values with the exact type specified.
1437   - If you don't know a value, use an empty array [] or null where appropriate,
1438   never change the type.
1439   - Constraints per field (types in parentheses):
1440     - required_capabilities (string[]): flat list of atomic capabilities.
1441     - matching_tools (string[]): names of existing tools.
1442     - missing_tools (object[]): each item MUST be an object with fields:
1443       - name (string)
1444       - description (string)
1445       - reason (string)
1446       - example_input_output (string) - a single-line string formatted as: "Input:
1447     {...} -> Output: {...}".
1448     - justification (string)
1449     - need_creation (boolean) - Phase 1 only: whether new tools are required.
1450
1451   Output Format
1452
1453   {
1454     "need_creation": true | false,
1455     "required_capabilities": ["verb_object", "..."],
1456     "matching_tools": ["existing_tool_a", "existing_tool_b"],
1457     "missing_tools": [
1458       {
1459         "name": "tool_for_gap",
1460         "description": "Does exactly this: ...",
1461         "reason": "Essential because: [Why Code? principle justification]",
1462         "example_input_output": "Input: {\\"x\\": 1} -> Output: {\\"y\\": 2}"
1463       }
1464     ],
1465     "justification": "1-3 sentences explaining overall coverage and gaps."
1466   }

```

1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511

#### Validation Checklist

- `matching_tools` in current available tool names (exact match).
- `missing_tools` names are unique and non-empty; reasons reference the Why Code? principles.
- If existing tools fully cover all `required_capabilities`: `need_creation=false` and `missing_tools=[]`.
- Output must be a single JSON object with no extra commentary.
- Each tool does exactly one specific operation.

#### EXAMPLES OF GOOD TOOLS:

- `count_even_numbers` - counts even numbers in a list`
- `extract_text_from_image` - extracts text from image (uses built-in vision)`
- `calculate_percentage` - calculates percentage from two numbers`
- `parse_addresses` - extracts and processes address data`

#### EXAMPLES OF BAD TOOLS:

- `address_parser` + `parity_checker` + `sunset_awning_counter` (redundant!)`
- `complex_data_analyzer` (too vague)`
- `multi_purpose_processor` (does too many things)`

Task:

```
{task}
```

Current Available Tools:

```
{available_tools}
```

### Prompt for Tool Spec Creation

[Back to ToC](#)

You are a Principal Engineer responsible for writing implementable specifications for new tools. Use the brief descriptions to generate detailed creation specs for ONLY those tools.

Design Principles (reused and specialized)

- One spec per missing tool; keep names consistent with Phase 1.
- Avoid overlapping with Current Available Tools; if overlap is detected, omit that spec.

STRICT TYPE CONTRACT

- For every field listed below, emit values with the exact type specified.
- If you don't know a value, use an empty array [] or null where appropriate, never change the type.

Output Format:

```
{
  "tool_creation_specs": [
    {
      "tool_name": "tool_for_gap",
      "tool_description": "1-2 sentence summary.",
      "input_parameters": [
        { "name": "param", "type": "string", "description": "...", "required":
true }
      ],
      "output_format": { "type": "object", "properties": { "result": { "type": "
number" } } },
      "core_logic": [
        "Step 1: Input validation ...",
        "Step 2: Core processing ..."
      ]
    }
  ]
}
```

1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565

```
}
```

```
Task:
{task}
```

```
Current Available Tools:
{available_tools}
```

## Prompt for Tool Creation

[Back to ToC](#)

You are a senior software engineering specialist building robust, reusable MCP (Model Context Protocol) tools following Claude Desktop standards.

### OBJECTIVE

Create a powerful and high-quality MCP tool to accurately and effectively solve the tasks specified below. This tool should be a direct and effective solution to a given problem, and it should be generated strictly according to the **Core Implementation Logic**.

### TASK

```
{original_query}
```

### TOOL NAME

```
{tool_name}
```

### Tool Description:

```
{tool_description}
```

### Core Implementation Logic:

```
{core_logic}
```

### Input Parameters:

```
{input_parameters}
```

### Output Format:

```
{output_format}
```

### DESIGN PRINCIPLES

- Solve the Problem Precisely: The tool's core goal is to provide a direct and effective solution for the current task. Prioritize functional correctness and robustness over premature generalization.
- Clear Interface Design: Clearly define the tool's input parameters and output format. Even if the tool has a specific function, its interface should be clear, well-documented, and easy to use.
- Professional Naming Convention: Choose a descriptive, domain-prefixed snake\_case name that accurately reflects the tool's specific purpose.

### OUTPUT CONTRACT (STRICT)

Return EXACTLY ONE JSON object as plain text (no markdown fences/backticks or extra prose). The object must have these keys:

1. name: string - professional general-purpose snake\_case with domain prefix
  - Must be EXACTLY "{tool\_name}" (do not rename or vary)
2. description: string - DETAILED description (5-10 sentences) covering:
  - What the tool does (core functionality)
  - Key capabilities and features
  - Typical use cases and scenarios
  - Input/output data types and formats

```

1566
1567     - Any limitations or constraints
1568
1569 3. input_schema: object - JSON Schema (Claude Desktop standard) defining
1570 parameters:
1571     {{
1572     "type": "object",
1573     "properties": {{
1574     "param_name": {{
1575     "type": "string|number|boolean|array|object",
1576     "description": "Detailed parameter description",
1577     "enum": ["optional", "allowed", "values"],
1578     "default": "optional_default_value",
1579     "example": "input example"
1580     }}
1581     }},
1582     "required": ["list_of_required_params"]
1583     }}
1584
1585     - Include ALL parameters the tool accepts
1586     - Provide detailed descriptions for each parameter
1587     - Specify types, constraints, enums, and defaults
1588     - Mark required vs optional parameters clearly
1589
1590 4. returns: object - Output format specification:
1591     {{
1592     "type": "string|object",
1593     "description": "Detailed description of return value",
1594     "format": "json|text|structured",
1595     "schema": {"optional": "output schema for structured returns"}}
1596     }}
1597
1598 5. module_code: string - COMPLETE Python module source to be saved as `<tool_name>.py`
1599
1600 Module requirements (STRICT):
1601 - Include necessary imports.
1602 - Define ONE public function implementing the tool with an EXPLICIT parameter list
1603 derived from input_schema (no **kwargs).
1604 - Implement robust validation and error handling per input_schema (types, required
1605 fields, enums, ranges).
1606 - Provide clear control flow and helpful error messages referencing parameter
1607 names.
1608 - Define TOOL_CONFIG = {{ "name": "{tool_name}", "description": <desc>, "function":
1609 <function_object>, "input_schema": <schema>, "returns": <returns> }} at module
1610 scope.
1611 - Multi-line Python code with normal 4-space indentation.
1612 - No external network calls or dependencies unless they are clearly documented and
1613 optional.
1614
1615 Naming constraints (STRICT):
1616 - The JSON field name MUST equal "{tool_name}".
1617 - In the Python module, TOOL_CONFIG['name'] MUST equal "{tool_name}".
1618 - Assume the file will be saved as "{tool_name}.py"; do not reference any other
1619 module name.
1620
1621 Formatting constraints:
1622 - No markdown fences.
1623 - No trailing backslashes at line ends.
1624 - The module must be self-contained and importable.
1625
1626 CAPABILITIES & QUALITY
1627 - Robust Input Handling: The tool must robustly handle various expected inputs and
1628 provide clear, helpful error messages for invalid or edge-case inputs.
1629

```

1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

- Graceful Degradation: Provide graceful degradation when inputs or external resources (if any) are missing or invalid.
- Strict Validation: Validate all inputs against the `input_schema` at the beginning of the function.
- Code Excellence: Optimize for clarity, maintainability, and performance; follow PEP 8 standards.
- Helpful Errors: Add helpful error messages that reference parameter names to guide the user.

#### REMINDERS

- Output must be a single JSON object (no additional commentary).
- Absolutely no json or python fences in the output.
- Ensure `input_schema` follows JSON Schema specification exactly.
- Provide detailed, production-ready documentation in all fields.

## Prompt for Agent Creation

[Back to ToC](#)

You are an expert agent designer. Create a minimal, reusable specialist agent specification for the task.

#### Inputs

- User Task: `{sub_task}`
- Available Tools: `{tools}`

#### Output format

- Return exactly one JSON object as the entire response.
- The object must contain only these four keys, in this exact order: `role`, `suggestions`, `tools`, `expertise`.
- Explicit schema:
  - `role`: string
  - `suggestions`: array of 3-5 strings
  - `tools`: array of 0-3 strings (tool names from Available Tools)
  - `expertise`: string (1-2 sentences; concise domain strengths and typical methodology)
- Do not include any other text, comments, code fences, or fields.

#### Constraints

##### Role

- A clear, human-readable professional title suitable for a business card.
- 2-5 words, Title Case, English only, no emojis/symbols.

##### Suggestions

- An array of 3-5 short, concrete execution hints for approaching this class of tasks.
- Each item starts with an imperative verb and is 6-16 words.
- Specific within the domain implied by `{sub_task}`, yet general-purpose (not tailored to a single query).
- No placeholders, no meta references (e.g., "JSON", "this prompt"), avoid tool names unless essential to the method.

##### Tools

- Choose a minimal subset of tool names taken verbatim from Available Tools.
- Use exact casing/spelling; do not invent, modify, or describe tools; no duplicates.
- If no tool is applicable or Available Tools is empty, use an empty array `[]`.

##### Expertise

- 1-2 sentences (20-160 characters) summarizing domain strengths and typical methodology.

1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727

- Mention information gathering, evaluation, analysis (quantitative/qualitative as applicable), and synthesis.
- Avoid listing tools; focus on capabilities and working approach.

#### General

- Ensure the role, suggestions, tools, and expertise are coherent with one another and with `{sub_task}`.
- Do not echo the inputs. Do not add explanations.
- Output must be valid JSON using double quotes and no trailing commas.

### Prompt for React Template

[Back to ToC](#)

You are a `{role}`. Based on prior agents' results and completed steps, your goal is to complete the current task as effectively as possible.

# Suggestions  
`{suggestions}`

# Progress So Far  
`{previous}`

# Past Experience References  
Below are relevant experiences from previous similar tasks. Reference the success experiences to learn effective approaches, and learn from the failure experiences to avoid common pitfalls.

## Success Experiences  
`{success_experiences}`

## Failure Experiences  
`{failure_experiences}`

You have access to the following tools:

# Tools  
`{tool}`

- # Steps
1. Review the outputs from previous agents to understand progress and context.
  2. Analyze the task and decompose it if needed. Utilize the available tools as appropriate.
  3. Define the current step you will complete, labeling it as 'CurrentStep'.
  4. Choose one Action from the available tools to execute the current step.

# Format example  
`{format_example}`

### Prompt for LLM as a Judge

[Back to ToC](#)

You are a Judge LLM responsible for evaluating the entire execution trajectory of a task and determining whether it was completed correctly.

#### OBJECTIVE

Analyze the complete task execution trajectory and provide a comprehensive evaluation of task completion, agent performance, and tool effectiveness.

#### TRAJECTORY INFORMATION

```

1728
1729 Original Query:
1730 {query}
1731
1732 Decomposed Subtasks:
1733 {subtasks}
1734
1735 Agent Assignments:
1736 {agent_assignments}
1737
1738 Agent Execution Processes:
1739 {agent_processes}
1740
1741 Result Aggregation Process:
1742 {aggregation_process}
1743
1744 Final Result:
1745 {final_result}
1746
1747 Newly Created Tools (if any):
1748 {created_tools}
1749
1750 EVALUATION CRITERIA
1751
1752 1. Task Completion Assessment:
1753 - Does the final result correctly answer the original query?
1754 - Are all decomposed subtasks properly addressed?
1755 - Is the aggregated result logically coherent and complete?
1756
1757 2. Agent Performance Assessment:
1758 - Did each agent execute its assigned subtask correctly?
1759 - Were the agents' reasoning processes sound and effective?
1760 - Did agents properly utilize available tools?
1761
1762 3. Tool Effectiveness Assessment:
1763 - Did existing tools function correctly when used?
1764 - Were newly created tools (if any) implemented correctly?
1765 - Did tools produce expected outputs?
1766
1767 EVALUATION OUTPUT FORMAT
1768
1769 You must provide your evaluation in the following JSON format:
1770
1771 {{
1772   "task_completed": true/false,
1773   "completion_quality": "good/poor",
1774   "overall_assessment": "Detailed overall assessment of task completion",
1775   "agent_evaluations": [
1776     {{
1777       "agent_id": "agent identifier",
1778       "agent_role": "role name",
1779       "subtask_id": "subtask identifier",
1780       "performance": "success/failure",
1781       "strengths": ["strength 1", "strength 2", ...],
1782       "issues": ["issue 1", "issue 2", ...]
1783     }},
1784     ...
1785   ],
1786   "tool_evaluations": [
1787     {{
1788       "tool_name": "tool name",
1789       "tool_type": "existing/newly_created",
1790       "effectiveness": "success/partial_success/failure",
1791       "issues": ["issue 1", "issue 2", ...],

```

1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1830  
1831  
1832  
1833  
1834  
1835

```
    }},  
    ...  
  ]  
}}
```

#### IMPORTANT GUIDELINES

1. Be objective and thorough in your evaluation
2. Provide specific, actionable feedback
3. Identify both strengths and weaknesses
4. Focus on patterns that can inform future improvements
5. Distinguish between agent failures and tool failures
6. Consider the complexity of the task when evaluating performance
7. Ensure all JSON is properly formatted and valid

Return ONLY the JSON evaluation, no additional text.

### Prompt for Tool Memory Generation

[Back to ToC](#)

You are a technical documentation specialist creating comprehensive tool implementation guides.

#### OBJECTIVE

Generate a detailed markdown section that captures the implementation experience of this tool for future reference and reuse. This will be one entry within a broader topic file.

#### TOOL INFORMATION

Tool Name: {tool\_name}  
Tool Description: {tool\_description}  
Required Capabilities: {required\_capabilities}  
Usage Examples: {usage\_examples}

Tool Implementation Code:  
{tool\_code}

#### DOCUMENT STRUCTURE

Create a markdown section with the following structure:

```
## [Generate a specific "How to..." question title here]
```

The title should:

- Start with "How to" or "How can I"
- Be specific to what this exact implementation does
- Be 5-15 words long
- Clearly distinguish this implementation from similar ones
- Example: "How to Download Files From a URL?", "How to Parse CSV Files with Custom Delimiters?"

#### ### Description

A concise 2-3 sentence overview explaining what this specific implementation does and what problem it solves.

#### ### Use Cases

A bulleted list of practical scenarios where this implementation is applicable.

Include:

- Specific real-world situations
- Types of data or inputs it handles
- Common integration patterns

```

1836
1837   ### Code Implementation
1838   The core implementation of the tool. Format as:
1839   ```python
1840   [Include the complete, clean tool implementation code here]
1841   ```
1842
1843   ### Tool Configuration
1844   (Optional - only include if the tool uses API keys, URLs, or other configuration)
1845   Document any configuration requirements:
1846   - API keys needed
1847   - Environment variables
1848   - URL endpoints
1849   - Configuration parameters
1850
1851   QUALITY REQUIREMENTS
1852   1. Write in clear, professional technical documentation style
1853   2. Use proper markdown formatting (## for main title, ### for subsections)
1854   3. Be concise but comprehensive
1855   4. Focus on reusability and understanding
1856   5. Include practical context, not just code
1857   6. Highlight key implementation patterns
1858   7. Note any important dependencies or requirements
1859
1860   OUTPUT FORMAT
1861   Return ONLY the complete markdown section. Do not include any preamble,
1862   explanations, or commentary outside the document itself. The section MUST start
1863   with "## How to..." as the first line.

```

## Prompt for Success Agent Memory Generation

[Back to ToC](#)

```

1864
1865   You are an experience synthesis specialist creating memory entries for an agent
1866   based on successful task execution.
1867
1868   OBJECTIVE
1869   Generate a comprehensive memory entry that captures the successful experience of
1870   this agent, including successful strategies, effective tool usage, and valuable
1871   insights for future similar tasks.
1872
1873   Agent Role: {agent_role}
1874   Agent Skills: {agent_skills}
1875
1876   TASK EXECUTION INFORMATION
1877   Subtask Description: {subtask_description}
1878   Task Context: {task_context}
1879   Tools Used: {tools_used}
1880   Execution Process:
1881   {execution_process}
1882
1883   Final Result: {final_result}
1884
1885   JUDGE EVALUATION
1886   Performance Rating: {performance_rating}
1887   Strengths Identified: {strengths}
1888   Key Success Patterns: {success_patterns}
1889
1890   MEMORY ENTRY STRUCTURE
1891
1892   Generate a memory entry in the following markdown format:
1893
1894   ## [Create a specific, descriptive title for this experience]

```

1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943

```
The title should:  
- Be specific to the type of task handled  
- Highlight the key capability demonstrated  
- Be 5-15 words long  
- Example: "How to Extract and Summarize Information from Multiple Web Sources"  
  
### Description  
A concise 2-3 sentence overview of what this experience teaches about handling  
this type of task.  
  
### Task Context  
Describe the original problem and the conditions under which this task was  
performed.  
  
### Experience of Success  
Detail the successful strategy and reasoning process:  
- Successful Strategy: Key decision points, problem-solving approach, and how  
tools were selected  
- Tools and Techniques: Tools used and how they contributed to success  
- Key Insights: What made this approach effective and when to apply similar  
strategies  
- Applicable Scenarios: Similar task types where this experience would be  
valuable  
- Potential Pitfalls: Conditions that might cause similar approaches to fail  
and warning signs to watch for  
  
QUALITY REQUIREMENTS  
1. Write in clear, professional documentation style  
2. Use proper markdown formatting  
3. Be specific and actionable  
4. Focus on transferable knowledge  
5. Highlight decision-making processes  
6. Include both what worked and why it worked  
  
OUTPUT FORMAT  
Return ONLY the complete markdown memory entry. Do not include any preamble or  
explanations outside the memory entry itself. The entry MUST start with "###"  
as the first line.
```

**Prompt for Failure Agent Memory Generation** [Back to ToC](#)

You are an error analysis specialist creating memory entries that help agents learn from failures.

**OBJECTIVE**  
Generate a comprehensive memory entry that captures the failure experience, analyzes root causes, and provides clear guidance on how to avoid similar failures in the future.

**AGENT INFORMATION**  
Agent Role: {agent\_role}  
Agent Skills: {agent\_skills}

**TASK EXECUTION INFORMATION**  
Subtask Description: {subtask\_description}  
Task Context: {task\_context}  
Tools Attempted: {tools\_used}  
Execution Process:  
{execution\_process}

1944  
 1945  
 1946 Failure Outcome: {failure\_outcome}  
 1947  
 1948 JUDGE EVALUATION  
 1949 Performance Rating: {performance\_rating}  
 1950 Issues Identified: {issues}  
 1951 Root Cause Analysis: {root\_cause}  
 1952 Suggested Fixes: {suggested\_fixes}

1953  
 1954 MEMORY ENTRY STRUCTURE  
 1955  
 1956 Generate a memory entry in the following markdown format:  
 1957  
 1958 ## [Create a specific, descriptive title for this failure case]  
 1959  
 1960 The title should:  
 1961 - Clearly indicate the problematic scenario  
 1962 - Be specific enough to match similar future situations  
 1963 - Be 5-15 words long  
 1964 - Example: "Common Pitfalls When Parsing Inconsistent CSV Data Formats"

1965  
 1966 ### Description  
 1967 A concise 2-3 sentence overview of what went wrong and the key lesson to take away.  
 1968  
 1969 ### Task Context  
 1970 Describe the task and conditions that led to this execution:  
 1971 - What was attempted  
 1972 - What was expected  
 1973 - What assumptions were made

1974  
 1975 ### Experience of Failure  
 1976 Detail what went wrong and how to avoid it:  
 1977 - **What Went Wrong**: Specific errors, incorrect decisions, tools misused, or  
 1978 logic errors  
 1979 - **Root Cause Analysis**: Why the approach didn't work, what was misunderstood,  
 1980 gaps in knowledge  
 1981 - **Corrective Actions**: Alternative approaches, additional checks, tools to use  
 1982 instead, validation steps  
 1983 - **Warning Signs**: Task characteristics that trigger this issue, context  
 1984 patterns, red flags during execution  
 1985 - **Related Success Patterns**: Successful approaches that should be used instead  
 1986 - **Partial Successes**: Parts of the approach that were correct, tools or  
 1987 techniques that functioned properly

1988  
 1989 QUALITY REQUIREMENTS  
 1990 1. Write in clear, instructive style  
 1991 2. Use proper markdown formatting  
 1992 3. Be honest and specific about failures  
 1993 4. Provide actionable corrective guidance  
 1994 5. Focus on learning and improvement  
 1995 6. Help prevent similar failures

1996  
 1997 OUTPUT FORMAT  
 Return ONLY the complete markdown memory entry. Do not include any preamble or  
 explanations outside the memory entry itself. The entry MUST start with "##" as  
 the first line.

## D CASE STUDY

### D.1 TOOL IMPLEMENTATION FOR *Simulate Piston Platform Game*

[BACK TO TOC](#)

As shown in Code 1, when performing complex probabilistic reasoning tasks, the agent extends its modeling capabilities by constructing a simulation tool and employing Monte Carlo methods to estimate the true underlying probabilities. The implementation of this tool is entirely grounded in the original problem formulation and faithfully simulates all rules specified in the task. Moreover, the implementation exhibits strong robustness, supporting simulations with arbitrary numbers of trials and arbitrary numbers of balls, and is, to some extent, adaptable to multiple variants of the original problem.

```

2009 1 import random
2010 2
2011 3 def simulate_piston_platform_game(num_balls=100, num_simulations=100000):
2012 4     """
2013 5     Simulates a specific ping-pong game mechanism involving a ball queue,
2014 6     a limited-capacity platform,
2015 7     and three random pistons with complex ejection/replacement rules.
2016 8     Used to calculate the
2017 9     probability of each ball number being 'ejected by a piston' (winning)
2018 10    versus being 'released' (eliminated).
2019 11
2020 12    Args:
2021 13    num_balls (int): Total number of balls in the queue, numbered 1
2022 14    to N.
2023 15    num_simulations (int): Monte Carlo iterations.
2024 16
2025 17    Returns:
2026 18    dict: A dictionary containing 'win_probabilities' and '
2027 19    best_choice'.
2028 20    """
2029 21
2030 22    # Step 1: Validate inputs
2031 23    if not isinstance(num_balls, int) or num_balls < 1:
2032 24        raise ValueError(f"Input 'num_balls' must be a positive integer.
2033 25        Got: {num_balls}")
2034 26
2035 27    if not isinstance(num_simulations, int) or num_simulations < 1:
2036 28        raise ValueError(f"Input 'num_simulations' must be an integer >=
2037 29    1. Got: {num_simulations}")
2038 30
2039 31    # Step 2: Initialize win_counts
2040 32    # Keys are integers 1 to N, initialized to 0
2041 33    win_counts = {i: 0 for i in range(1, num_balls + 1)}
2042 34
2043 35    # Step 3: Begin simulation loop
2044 36    for _ in range(num_simulations):
2045 37        # Step 3a: Initialize queue and platform
2046 38        deck = list(range(1, num_balls + 1))
2047 39        platform = []
2048 40
2049 41        # Initial fill of the platform (max 3 balls)
2050 42        for _ in range(3):
2051 43            if deck:
2052 44                platform.append(deck.pop(0))
2053 45            else:
2054 46                platform.append(None)
2055 47
2056 48        # Helper function to refill balls from deck
2057 49        def get_next_balls(count):
2058 50            new_balls = []
2059 51            for _ in range(count):
2060 52                if deck:
2061 53                    new_balls.append(deck.pop(0))

```

```

2052 46         else:
2053 47             new_balls.append(None)
2054 48         return new_balls
2055 49
2056 50     # Step 4: Simulate piston firing process
2057 51     # Continue while there are still balls on the platform
2058 52     while any(ball is not None for ball in platform):
2059 53         # Select a random piston (1, 2, or 3)
2060 54         piston = random.randint(1, 3)
2061 55
2062 56         # Current platform state
2063 57         p1, p2, p3 = platform[0], platform[1], platform[2]
2064 58
2065 59         # Step 5: Apply Piston Rules
2066 60         if piston == 1:
2067 61             # If Piston 1: Eject Pos 1 (Win). Move Pos 2->1, Pos
2068 62             3->2. Refill 1 ball.
2069 63             if p1 is not None:
2070 64                 win_counts[p1] += 1
2071 65
2072 66                 incoming = get_next_balls(1)
2073 67                 platform = [p2, p3, incoming[0]]
2074 68
2075 69         elif piston == 2:
2076 70             # If Piston 2: Eject Pos 2 (Win). Release Pos 1 (Loss).
2077 71             Move Pos 3->1. Refill 2 balls.
2078 72             if p2 is not None:
2079 73                 win_counts[p2] += 1
2080 74
2081 75                 # p1 is released (eliminated), so we don't increment its
2082 76                 win count
2083 77
2084 78                 incoming = get_next_balls(2)
2085 79                 platform = [p3, incoming[0], incoming[1]]
2086 80
2087 81         elif piston == 3:
2088 82             # If Piston 3: Eject Pos 3 (Win). Release Pos 1 (Loss).
2089 83             Move Pos 2->1. Refill 2 balls.
2090 84             if p3 is not None:
2091 85                 win_counts[p3] += 1
2092 86
2093 87                 # p1 is released (eliminated), so we don't increment its
2094 88                 win count
2095 89
2096 90                 incoming = get_next_balls(2)
2097 91                 platform = [p2, incoming[0], incoming[1]]
2098 92
2099 93     # Step 8: Calculate probabilities
2100 94     win_probabilities = {}
2101 95     for i in range(1, num_balls + 1):
2102 96         win_probabilities[str(i)] = win_counts[i] / num_simulations
2103 97
2104 98     # Find the ball with the highest probability
2105 99     best_ball_num = max(win_counts, key=win_counts.get)
2106 100
2107 101     # Step 9: Format output object
2108 102     return {
2109 103         "win_probabilities": win_probabilities,
2110 104         "best_choice": best_ball_num
2111 105     }

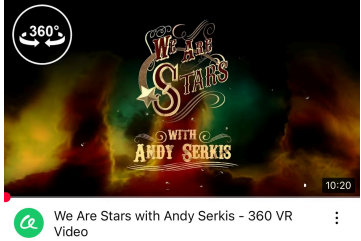
```

Code 1: Code Implenmenation for Simulate Piston Platform Game Tool


D.2 TOOL IMPLEMENTATION FOR *Youtube Audio Transcriber*[BACK TO TOC](#)

2106  
2107  
2108  
2109  
2110  
2111  
2112  
2113  
2114  
2115  
2116  
2117  
2118  
2119  
2120  
2121  
2122  
2123  
2124  
2125  
2126  
2127  
2128  
2129  
2130  
2131  
2132  
2133  
2134  
2135  
2136  
2137  
2138  
2139  
2140  
2141  
2142  
2143  
2144  
2145  
2146  
2147  
2148  
2149  
2150  
2151  
2152  
2153  
2154  
2155  
2156  
2157  
2158  
2159


**Question:** In the YouTube 360 VR video from March 2018 narrated by the voice actor of Lord of the Rings' Gollum, what number was mentioned by the narrator directly after dinosaurs were first shown in the video?  
**Ground Truth:** 10000000





---

 `web_search (query="gollum voice actor")`


**Tool Response:** [Search Information]


 I need to watch "We Are Stars with Andy Serkis -360 VR Video" to determine the time of the first appearance of dinosaurs.`web_search(query="We Are Stars with Andy Serkis -360 VR Video")`

 I cannot access the video, but based on common sense, the first number that appears should be 0.


✘ **Incorrect!**

(a) experience-centric frameworks

 To accurately identify the number mentioned after the dinosaur scene, I need to analyze the video's full content. **I will create a tool to transcribe the YouTube audio and extract the text.**

 `youtube_audio_transcriber (url=youtube_url)`  
*Full Implementation in Code 2*

**Tool Response:** {"video\_metadata": ..., "transcription\_text": ....., "file\_info": .....,}

 Based on the transcription text, I believe the answer is 10000000.

✔ **Correct!**

(b) Mem<sup>2</sup>Evolve

Figure 7: **Case Study 2 on YouTube Video Subtitle Extraction.** When initialized with only a web search tool, (a) experience-centric frameworks fail to handle tasks situated beyond their capability boundary, such as retrieving internal video content, leading to incorrect answers based on general common sense. In contrast, (b) Mem<sup>2</sup>Evolve leverages the guidance of accumulated experience to dynamically generate high-quality tools (e.g., a custom subtitle transcriber), effectively breaking through capability boundaries to access the necessary context and derive the correct answer.

As illustrated in Code 2, when the model encounters complex tasks in open-world scenarios, such as "In the YouTube 360 VR video narrated by the voice actor of Gollum from *The Lord of the Rings* in March 2018, what number does the narrator directly mention after the dinosaur first appears?", it is necessary to perform semantic analysis on the audio content of YouTube videos to obtain the correct answer. To address such problems, the Agent autonomously generates the *Youtube Audio Transcriber* tool, which integrates open-source libraries including yt-dlp, pytube, and speech\_recognition to enable audio transcription for arbitrary YouTube video URLs with multi-language support, thereby generalizing to similar tasks. The design of this tool demonstrates three key characteristics: First, the implementation leverages mature open-source components rather than building from scratch, significantly improving development efficiency and reliability. Second, the tool exhibits strong robustness, effectively handling diverse input formats and various edge cases. Third, the tool adheres to the principle of generality, targeting a category of tasks rather than a single problem, thus possessing good reusability to support the execution of subsequent similar tasks.

```

1 import os
2 import math
3 import shutil
4 import json

```

```

2160 5 from urllib.parse import urlparse, parse_qs
2161 6 import yt_dlp
2162 7 from pydub import AudioSegment
2163 8 import speech_recognition as sr
2164 9
2165 10
2166 11 def youtube_audio_transcriber(youtube_url: str, language_code: str = "en-
2167 12 US", chunk_duration_sec: int = 15):
2168 13     """
2169 14     Downloads audio from YouTube and transcribes it into text.
2170 15     Args:
2171 16         youtube_url (str): The full URL of the YouTube video.
2172 17         language_code (str): Language code for recognition (default 'en-
2173 18 US').
2174 19         chunk_duration_sec (int): Duration of each processing chunk in
2175 20 seconds (default 15).
2176 21     Returns:
2177 22         dict: Object containing transcription_text, video_metadata, and
2178 23 file_info.
2179 24     """
2180 25     workspace_dir = "workspace"
2181 26     temp_dir = "temp_chunks"
2182 27
2183 28     # Ensure directories exist
2184 29     if not os.path.exists(workspace_dir):
2185 30         os.makedirs(workspace_dir)
2186 31
2187 32     video_id = None
2188 33     downloaded_mp3_path = None
2189 34     full_transcription = ""
2190 35     video_metadata = {}
2191 36
2192 37     try:
2193 38         # --- Step 1: Validate inputs and parse the `youtube_url` to
2194 39         # extract the Video ID ---
2195 40         if 'youtube.com' in youtube_url:
2196 41             parsed_url = urlparse(youtube_url)
2197 42             video_id = parse_qs(parsed_url.query).get('v', [None])[0]
2198 43         elif 'youtu.be' in youtube_url:
2199 44             parsed_url = urlparse(youtube_url)
2200 45             video_id = parsed_url.path.lstrip('/')
2201 46
2202 47         if not video_id:
2203 48             raise ValueError(f"Could not extract Video ID from URL: {
2204 49 youtube_url}")
2205 50
2206 51         # --- Step 2: Configure `yt_dlp` options ---
2207 52         # Set output template to workspace, convert to mp3 192kbps
2208 53         output_template = os.path.join(workspace_dir, f"{video_id}_%(
2209 54 title)s.%(ext)s")
2210 55
2211 56         ydl_opts = {
2212 57             'format': 'bestaudio/best',
2213 58             'outtmpl': output_template,
2214 59             'postprocessors': [{
2215 60                 'key': 'FFmpegExtractAudio',
2216 61                 'preferredcodec': 'mp3',
2217 62                 'preferredquality': '192',
2218 63             }],
2219 64             'quiet': True,
2220 65             'no_warnings': True,
2221 66         }

```

```

2214 63
2215 64 # --- Step 3: Execute the download & Extract Metadata ---
2216 65 print(f"[Tool] Starting download for Video ID: {video_id}...")
2217 66 with yt_dlp.YoutubeDL(ydl_opts) as ydl:
2218 67     info = ydl.extract_info(youtube_url, download=True)
2219 68
2220 69     # Format duration string
2221 70     duration = info.get('duration', 0)
2222 71     m, s = divmod(duration, 60)
2223 72     h, m = divmod(m, 60)
2224 73     duration_str = f"{h:02d}:{m:02d}:{s:02d}" if h else f"{m:02d}
2225 74     }:{s:02d}"
2226 75
2227 76     video_metadata = {
2228 77         "title": info.get('title', 'Unknown'),
2229 78         "uploader": info.get('uploader', 'Unknown'),
2230 79         "duration_str": duration_str,
2231 80         "video_id": video_id
2232 81     }
2233 82
2234 83     # --- Step 4: Locate the downloaded MP3 file ---
2235 84     # yt-dlp might replace characters in the filename, so we search
2236 85     by video_id
2237 86     for file in os.listdir(workspace_dir):
2238 87         if video_id in file and file.endswith('.mp3'):
2239 88             downloaded_mp3_path = os.path.join(workspace_dir, file)
2240 89             break
2241 90
2242 91     if not downloaded_mp3_path:
2243 92         raise FileNotFoundError("Audio file not found after download
2244 93         process.")
2245 94
2246 95     # --- Step 5: Initialize SpeechRecognition & Temp Directory ---
2247 96     recognizer = sr.Recognizer()
2248 97
2249 98     if os.path.exists(temp_dir):
2250 99         shutil.rmtree(temp_dir)
2251 100     os.makedirs(temp_dir)
2252 101
2253 102     # --- Step 6: Load MP3 & Calculate Chunks ---
2254 103     print("[Tool] Loading audio file for processing...")
2255 104     audio = AudioSegment.from_file(downloaded_mp3_path)
2256 105
2257 106     # pydub works in milliseconds
2258 107     chunk_length_ms = chunk_duration_sec * 1000
2259 108     num_chunks = math.ceil(len(audio) / chunk_length_ms)
2260 109
2261 110     print(f"[Tool] Audio length: {len(audio)/1000:.2f}s. Split into {
2262 111     num_chunks} chunks.")
2263 112
2264 113     # --- Step 7, 8, 9: Iterate, Slice, Recognize, Append ---
2265 114     print("[Tool] Starting transcription...")
2266 115
2267 116     for i in range(num_chunks):
2268 117         start_ms = i * chunk_length_ms
2269 118         end_ms = (i + 1) * chunk_length_ms
2270 119
2271 120         # Slice audio
2272 121         chunk = audio[start_ms:end_ms]
2273 122
2274 123         # Export to WAV (required by SpeechRecognition)
2275 124         chunk_filename = os.path.join(temp_dir, f"chunk_{i}.wav")
2276 125         chunk.export(chunk_filename, format="wav")
2277 126
2278 127         # Recognize

```

```

2268 124         with sr.AudioFile(chunk_filename) as source:
2269 125             audio_data = recognizer.record(source)
2270 126             try:
2271 127                 text = recognizer.recognize_google(audio_data,
2272 128                 language=language_code)
2273 129                 full_transcription += text + " "
2274 130             except sr.UnknownValueError:
2275 131                 # Audio was not understood (silence, noise, music)
2276 132                 pass
2277 133             except sr.RequestError as e:
2278 134                 print(f"[Tool] API Error on chunk {i}: {e}")
2279 135             except Exception as e:
2280 136                 print(f"[Tool] Unexpected error on chunk {i}: {e}")
2281 137
2282 138         full_transcription = full_transcription.strip()
2283 139
2284 140         # Get file size for report
2285 141         file_size_mb = os.path.getsize(downloaded_mp3_path) / (1024 *
2286 142         1024)
2287 143
2288 144         # --- Step 10: Clean up temporary files ---
2289 145         if os.path.exists(temp_dir):
2290 146             shutil.rmtree(temp_dir)
2291 147
2292 148         # --- Step 11: Return result object ---
2293 149         return {
2294 150             "transcription_text": full_transcription,
2295 151             "video_metadata": video_metadata,
2296 152             "file_info": {
2297 153                 "local_path": downloaded_mp3_path,
2298 154                 "file_size_mb": round(file_size_mb, 2)
2299 155             }
2300 156         }
2301 157     except Exception as e:
2302 158         # Cleanup temp if error occurs
2303 159         if os.path.exists(temp_dir):
2304 160             shutil.rmtree(temp_dir)
2305 161
2306 162         # Return error structure or raise
2307 163         return {
2308 164             "error": str(e),
2309 165             "transcription_text": "",
2310 166             "video_metadata": video_metadata if video_metadata else {},
2311 167             "file_info": {}
2312 168         }

```

Code 2: Tool Implementation for Youtube Audio Transcriber

### D.3 EXPERIENCE GUIDANCE TOOL CREATION

[BACK TO TOC](#)

In this case study, we demonstrate how Mem<sup>2</sup>Evolve creates a new tool under the guidance of Experience Memory. When the system already contains a memory item titled “*How to Analyze Images Using the GPT-4o Multimodal Model?*” (Figure 8), the model, in generating a tool (Code 3) for analyzing YouTube video content, first extracts screenshots from the video at fixed frame intervals. Subsequently, it applies GPT-4o to analyze each extracted frame and automatically constructs a complete, well-aligned prompt tailored to the analysis task.

```

2319 1 import os
2320 2 import json
2321 3 import subprocess
2322 4 import base64

```

```

2322
2323
2324 ## How to Analyze Images Using GPT-4o Multimodal Model?
2325
2326 ### Description
2327 Parse and analyze an image file using GPT-4o multimodal model. This code can understand
2328 complex visual content, generate captions, extract tables as HTML, create SVG code
2329 for geometric shapes, and answer specific questions about images.
2330
2331 ### Use Cases
2332 - Product image analysis for e-commerce catalog management
2333 - Medical image interpretation and diagnostic support
2334 - Security and surveillance image analysis
2335 - Educational content creation from visual materials
2336 - Art and cultural artifact documentation
2337 - Scientific image analysis and research documentation
2338 - Social media content moderation and analysis
2339
2340 ### Tool Implementation
2341 ```python
2342 # Partial code implementation is omitted here
2343
2344 # Prepare API request payload
2345 payload = {
2346     "model": "gpt-4o-2024-11-20",
2347     "messages": [
2348         {
2349             "role": "user",
2350             "content": [
2351                 {
2352                     "type": "text",
2353                     "text": prompt,
2354                 },
2355                 {
2356                     "type": "image_url",
2357                     "image_url": {
2358                         "url": f"{img_type}{img_base64}"
2359                     }
2360                 }
2361             ],
2362         },
2363     ],
2364     "max_tokens": 16384,
2365 }
2366
2367 # Get API credentials from environment variables
2368 api_key = os.getenv("OPENAI_API_KEY")
2369 api_base = os.getenv("OPENAI_BASE_URL")
2370
2371 headers = {
2372     "Content-Type": "application/json",
2373     "Authorization": f"Bearer {api_key}"
2374 }
2375
2376 # Send request to OpenAI API
2377 response = requests.post(f"{api_base}/chat/completions", headers=headers, json=payload)
2378
2379 result = response.json()
2380 output = result["choices"][0]["message"]["content"]
2381 ```
2382
2383
2384
2385

```

Figure 8: **Tool Experience: Using the GPT-4o for Image Analysis.** This tool experience illustrates how to call the GPT-4o API to analyze images, where the agent can customize prompts to steer GPT-4o toward diverse and complex visual understanding tasks (e.g., recognition, counting, spatial reasoning, chart/diagram interpretation, and multimodal grounding). Each tool experience is organized into four fields: *Title*, *Description*, *Use Cases*, and *Tool Implementation*.

```

2376 5 import requests
2377 6 import math
2378 7 import shutil
2379 8 import tempfile
2380 9 from dotenv import load_dotenv
2381 10
2382 11 load_dotenv()
2383 12
2384 13
2384 14 def analyze_video_for_species(youtube_url: str, target_subject: str = "
2385 15     bird species", sampling_interval: int = 10):
2386 16     """
2387 17     Analyzes a YouTube video to determine the maximum number of distinct
2388 18     species visible simultaneously.
2389 19     """
2390 20     api_key = os.getenv("OPENAI_API_KEY")
2391 21     if not api_key:
2392 22         return {"error": "Missing OPENAI_API_KEY environment variable."}
2393 23
2394 24     workspace_dir = "workspace"
2395 25     if not os.path.exists(workspace_dir):
2396 26         os.makedirs(workspace_dir)
2397 27
2398 28     max_species_count = 0
2399 29     best_frame_data = {
2400 30         "count": 0,
2401 31         "species_list": [],
2402 32         "description": "No data found."
2403 33     }
2404 34     best_timestamp = "00:00:00"
2405 35
2406 36     try:
2407 37         # --- Step 1: Validate URL and Extract Metadata ---
2408 38         print(f"[Tool] Getting video info for: {youtube_url}")
2409 39         cmd_info = ['yt-dlp', '--dump-json', '--no-playlist', youtube_url
2410 40     ]
2411 41         result_info = subprocess.run(cmd_info, capture_output=True, text=
2412 42     True, timeout=30)
2413 43
2414 44         if result_info.returncode != 0:
2415 45             raise ValueError(f"Failed to extract video info: {result_info
2416 46     .stderr}")
2417 47
2418 48         video_info = json.loads(result_info.stdout)
2419 49         duration = video_info.get('duration') # seconds
2420 50         video_id = video_info.get('id', 'unknown')
2421 51
2422 52         if not duration:
2423 53             raise ValueError("Could not determine video duration.")
2424 54
2425 55         # --- Step 2: Calculate Timestamps ---
2426 56         # Limit total checks to avoid excessive API usage cost in this
2427 57     demo implementation
2428 58         # For production, you might want to remove the limit or increase
2429 59     interval
2430 60         timestamps_sec = range(0, int(duration), sampling_interval)
2431 61         total_steps = len(timestamps_sec)
2432 62
2433 63         print(f"[Tool] Video Duration: {duration}s. Sampling every {
2434 64     sampling_interval}s. Total checks: {total_steps}")
2435 65
2436 66         # --- Step 3 & 4: Iterate through timestamps ---
2437 67         for idx, current_sec in enumerate(timestamps_sec):

```

```

2430         # Format timestamp HH:MM:SS
2431         m, s = divmod(current_sec, 60)
2432         h, m = divmod(m, 60)
2433         timestamp_str = f"{h:02d}:{m:02d}:{s:02d}"
2434
2435         print(f"[Tool] ({idx+1}/{total_steps}) Processing timestamp:
2436         {timestamp_str}...")
2437
2438         # --- Step 5: Download Segment (using ffmpeg downloader for
2439         speed) ---
2440         # Create a unique temp file for this segment
2441         temp_segment_path = os.path.join(workspace_dir, f"temp_{
2442         video_id}_{current_sec}.mp4")
2443         screenshot_path = os.path.join(workspace_dir, f"frame_{
2444         video_id}_{current_sec}.jpg")
2445
2446         try:
2447             # Use yt-dlp with ffmpeg external downloader to fetch
2448             just a tiny snippet
2449             # This avoids downloading the whole video
2450             download_cmd = [
2451                 'yt-dlp',
2452                 '--format', 'best[height<=720]', # 720p is enough for
2453                 recognition
2454                 '--external-downloader', 'ffmpeg',
2455                 '--external-downloader-args', f'ffmpeg_i:-ss {
2456         current_sec} -t 2', # download 2 seconds
2457                 '--output', temp_segment_path,
2458                 '--quiet', '--no-warnings',
2459                 youtube_url
2460             ]
2461             subprocess.run(download_cmd, capture_output=True, timeout
2462             =60)
2463
2464             # --- Step 6: Extract Screenshot ---
2465             # Check if video segment exists (sometimes yt-dlp appends
2466             ext)
2467             found_video = None
2468             for ext in ['.mp4', '.webm', '.mkv']:
2469                 check_path = temp_segment_path.replace('.mp4', ext) #
2470                 naive replacement
2471                 if os.path.exists(check_path): # Check exact match
2472                 first if template wasn't dynamic
2473                     found_video = check_path
2474                     break
2475                 # Handle yt-dlp output template behavior if needed
2476                 if os.path.exists(temp_segment_path):
2477                     found_video = temp_segment_path
2478                     break
2479
2480             if not found_video:
2481                 print(f"    -> Warning: Could not download segment
2482                 for {timestamp_str}, skipping.")
2483                 continue
2484
2485             # Capture first frame of the segment
2486             ffmpeg_cmd = [
2487                 'ffmpeg', '-i', found_video,
2488                 '-vframes', '1', '-q:v', '2', '-y',
2489                 screenshot_path
2490             ]
2491             subprocess.run(ffmpeg_cmd, capture_output=True, timeout
2492             =10)
2493
2494             if not os.path.exists(screenshot_path):

```

```

2484     print(f"    -> Warning: Screenshot failed for {
114 timestamp_str}, skipping.")
2485     continue
115
2486
116
2487
117
2488     # --- Step 7: Encode Base64 ---
118     with open(screenshot_path, "rb") as image_file:
2489         base64_image = base64.b64encode(image_file.read()).
119 decode('utf-8')
2490
120
2491
121
2492     # --- Step 8: Construct AI Prompt ---
122     prompt_text = (
123         f"Analyze this image specifically to count {
2493 target_subject}. "
124         f"Identify all distinct {target_subject} present. "
2494         "Ignore statues, pictures, or reflections if clearly
125 not real/live instances (unless the target is inanimate). "
2495         "Return ONLY valid JSON with the following keys: "
126         "'count' (integer, number of UNIQUE species visibly
2496 distinct simultaneously), "
127         "'species_list' (array of strings, names of
2497 identified species), "
128         "'description' (string, brief summary of the scene)."

```

```

2538 # --- Step 10: Compare and Update ---
2539 if current_count > max_species_count:
2540     max_species_count = current_count
2541     best_timestamp = timestamp_str
2542     best_frame_data = analysis
2543     print(f"    -> New Record! {max_species_count} at {
2544 best_timestamp}")
2545
2546     except Exception as e:
2547         print(f"    -> Error processing step {timestamp_str}: {e}
2548 ")
2549
2550     finally:
2551         # --- Step 11: Cleanup ---
2552         if os.path.exists(screenshot_path):
2553             os.remove(screenshot_path)
2554         if 'found_video' in locals() and found_video and os.path.
2555 exists(found_video):
2556             os.remove(found_video)
2557
2558     # --- Step 12: Return Results ---
2559     result = {
2560         "max_simultaneous_species": max_species_count,
2561         "best_timestamp": best_timestamp,
2562         "identified_species": best_frame_data.get("species_list", [])
2563     },
2564     "analysis_summary": best_frame_data.get("description", "")
2565 }
2566
2567     return result
2568
2569 except Exception as e:
2570     return {"error": str(e)}
2571
2572

```

Code 3: Code Implementation for Analyze Video for Species Tool

#### D.4 COMPARISON OF TOOL GENERATION WITH AND WITHOUT EXPERIENCE GUIDANCE

[BACK TO TOC](#)

This case provides an intuitive demonstration of the guiding role of Experience Memory in the generation of new tools. Specifically, as shown in Figure 9, Mem<sup>2</sup>Evo is required to read color-coded cells from an Excel file during task execution. With guidance from experience, the tool generated by Mem<sup>2</sup>Evo is able to accurately retrieve the target cells together with their original color information and output the results in a standardized format (full implementation in Code 4). In contrast, in the absence of such experiential guidance, the model generates a tool that relies solely on pandas to read the Excel content, failing to preserve and return the color attributes (full implementation in Code 5). This limitation ultimately leads to unsuccessful task execution.

```

2581 import os
2582 import pandas as pd
2583 from openpyxl import load_workbook
2584
2585 def parse_excel_with_styles(file_path: str, row_limit: int = 100):
2586     """
2587     Parses an Excel or CSV file and returns the content as formatted HTML
2588     with style information preserved.
2589     """
2590
2591     # Internal helper function to extract styles
2592     def get_cell_style(cell):
2593         """Extract style information from a cell and return as CSS style
2594         string."""

```


2592  
2593  
2594  
2595  
2596  
2597  
2598  
2599  
2600  
2601  
2602  
2603  
2604  
2605  
2606  
2607  
2608  
2609  
2610  
2611  
2612  
2613  
2614  
2615  
2616  
2617  
2618  
2619  
2620  
2621  
2622  
2623  
2624  
2625  
2626  
2627  
2628  
2629  
2630  
2631  
2632  
2633  
2634  
2635  
2636  
2637  
2638  
2639  
2640  
2641  
2642  
2643  
2644  
2645

**Question:** Each cell in the attached spreadsheet represents a plot of land. **The color of the cell indicates who owns that plot. Green cells are plots owned by Earl Smith.** Can Earl walk through every plot he owns (and no other plots) and return to his starting plot without backtracking? For this question, consider backtracking to be any instance where Earl would enter a plot of land he had already entered since leaving his starting plot.


**File Content:**


**Ground Truth: No**

---


 `read_excel_basic (file_path=file_path)`  
*Full Implementation in Code 4*

**Tool Response:** {'file\_name': '5cfb274c-0207-4aa7-9575-6ac0bd95d9b2.xlsx', 'sheet\_content': {'Sheet1': ''}}



 I don't have enough information!

 **Incorrect!**

(a) w/o experience-guidance

 `parse_excel_with_styles (file_path=file_path)`  
*Full Implementation in Code 3*

**Tool Response:** ..., <table border='1' style='border-collapse:collapse;'> <tr><th style='background-color:#4A86E8;'></th><th style='background-color:#4A86E8;'>>, .....

 Based on the information in the attachment, I believe the answer is No.  **Correct!**

(b) w/ experience-guidance

**Figure 9: Case Study 4: Experience-Guided Tool Generation for Attribute-Preserving Excel Parsing.** This case illustrates how Experience Memory guides Mem<sup>2</sup>Evolve to generate task-appropriate tools that preserve critical non-textual attributes. When required to extract color-coded cells from an Excel file, Mem<sup>2</sup>Evolve leverages past experience to synthesize a tool capable of accurately retrieving both cell values and their original color information in a standardized output format (full implementation in Code 4). In contrast, without experiential guidance, the generated tool relies solely on pandas, which fails to retain color attributes (full implementation in Code 5), leading to unsuccessful task execution.

```

2646 13         styles = []
2647 14
2648 15         # Check for bold formatting
2649 16         if cell.font and cell.font.bold:
2650 17             styles.append('font-weight:bold;')
2651 18
2652 19         # Check for italic formatting
2653 20         if cell.font and cell.font.italic:
2654 21             styles.append('font-style:italic;')
2655 22
2656 23         # Extract font color
2657 24         # Step 6 & 7: Handle Color Processing (ARGB -> RGB)
2658 25         color = getattr(cell.font, 'color', None)
2659 26         if color is not None and getattr(color, 'type', None) == 'rgb':
2660 27             rgb = getattr(color, 'rgb', None)
2661 28             if isinstance(rgb, str) and len(rgb) >= 6:
2662 29                 # Slice the last 6 characters to ignore Alpha channel (
2663 30                 ARGB -> RGB)
2664 31                 styles.append(f'color:#{rgb[-6:]};')
2665 32
2666 33         # Extract background color
2667 34         fill = getattr(cell, 'fill', None)
2668 35         fgColor = getattr(fill, 'fgColor', None)
2669 36         if fgColor is not None and getattr(fgColor, 'type', None) == 'rgb':
2670 37             rgb = getattr(fgColor, 'rgb', None)
2671 38             # Filter out transparent/invalid colors (00000000 usually
2672 39             means no fill in some contexts, but checking length is safer)
2673 40             if isinstance(rgb, str) and rgb != '00000000' and len(rgb) >=
2674 41             6:
2675 42                 styles.append(f'background-color:#{rgb[-6:]};')
2676 43
2677 44         return ''.join(styles)
2678 45
2679 46         # Step 1: Validate file existence and format
2680 47         if not os.path.exists(file_path):
2681 48             return {"error": f"Error: File '{file_path}' does not exist.", "
2682 49             html_content": "", "file_metadata": {}}
2683 50
2684 51         supported_formats = ['.xlsx', '.xls', '.csv']
2685 52         file_ext = os.path.splitext(file_path)[1].lower()
2686 53
2687 54         if file_ext not in supported_formats:
2688 55             return {"error": f"Error: Unsupported file format '{file_ext}'.",
2689 56             "html_content": "", "file_metadata": {}}
2690 57
2691 58         html_output = ""
2692 59         metadata = {
2693 60             "file_type": file_ext,
2694 61             "sheet_names": []
2695 62         }
2696 63
2697 64         try:
2698 65             # Step 2: Handle CSV files
2699 66             if file_ext == '.csv':
2700 67                 df = pd.read_csv(file_path)
2701 68                 metadata["sheet_names"] = ["csv_data"]
2702 69
2703 70                 html_output += f"<h2>CSV : {os.path.basename(file_path)}</h2>
2704 71                 <\n"
2705 72                 html_output += f"<p>Rows: {df.shape[0]}, Columns: {df.shape
2706 73                 [1]}</p><\n"
2707 74                 html_output += "<table border='1'><\n"
2708 75
2709 76             # Add header

```

```

2700         html_output += "<tr>"
2701     for col in df.columns:
2702         html_output += f"<th>{col}</th>"
2703     html_output += "</tr>\n"
2704
2705     # Add data rows
2706     for i, row in df.head(row_limit).iterrows():
2707         html_output += "<tr>"
2708         for value in row:
2709             val_str = str(value) if pd.notna(value) else ""
2710             html_output += f"<td>{val_str}</td>"
2711         html_output += "</tr>\n"
2712
2713     if len(df) > row_limit:
2714         html_output += f"<tr><td colspan='{len(df.columns)}'>...
2715         ({len(df) - row_limit} more rows)</td></tr>\n"
2716
2717     html_output += "</table>\n"
2718
2719     # Step 3: Handle Excel files
2720     else:
2721         # data_only=True is essential to get values instead of
2722         formulas
2723         wb = load_workbook(file_path, data_only=True)
2724         metadata["sheet_names"] = wb.sheetnames
2725
2726         html_output += f"<h1>Excel: {os.path.basename(file_path)}</h1
2727         >\n"
2728
2729         # Step 4: Iterate through sheets
2730         for sheet in wb.worksheets:
2731             html_output += f"<h2>Sheet: {sheet.title}</h2>\n"
2732
2733             max_row = sheet.max_row
2734             max_col = sheet.max_column
2735
2736             html_output += f"<p>Rows: {max_row}, Columns: {max_col}</
2737             p>\n"
2738             html_output += "<table border='1' style='border-collapse:
2739             collapse;'>\n"
2740
2741             # Step 5: Process rows and cells
2742             # enumerate(..., 1) makes i start at 1
2743             for i, row in enumerate(sheet.iter_rows(max_row=min(
2744             max_row, row_limit)), 1):
2745                 html_output += "<tr>"
2746                 for cell in row:
2747                     tag = "th" if i == 1 else "td" # Assume first
2748                     row is header
2749
2750                     # Step 6: Extract style
2751                     style = get_cell_style(cell)
2752                     value = cell.value if cell.value is not None else
2753                     ""
2754
2755                     # Step 8: Construct HTML with inline styles
2756                     if style:
2757                         html_output += f"<{tag} style='{style}'>{
2758                         value}</{tag}>"

```

Code 4: Tool Implementation for Parse Excel With Styles

```

2750
2751
2752
2753
1 import os

```

```

2754 2 import pandas as pd
2755 3
2756 4 def read_excel_basic(file_path: str, preview_rows: int = 50):
2757 5     """
2758 6     Basic Excel reader using standard Pandas functionality.
2759 7     Fails to capture style information required for color-based riddles.
2760 8     """
2761 9
2762 10    # Step 1: Validate file existence
2763 11    if not os.path.exists(file_path):
2764 12        return {"error": f"File '{file_path}' not found."}
2765 13
2766 14    file_ext = os.path.splitext(file_path)[1].lower()
2767 15    data_output = {}
2768 16
2769 17    try:
2770 18        # Step 2: Read file based on extension
2771 19        # Pandas read_excel defaults to reading values
2772 20        if file_ext == '.csv':
2773 21            df = pd.read_csv(file_path)
2774 22            # Convert to markdown-style string for readability
2775 23            data_output['csv_data'] = df.head(preview_rows).to_markdown(
2776 24                index=False)
2777 25
2778 26            elif file_ext in ['.xlsx', '.xls']:
2779 27                # sheet_name=None reads all sheets into a dictionary
2780 28                sheets = pd.read_excel(file_path, sheet_name=None)
2781 29
2782 30                for sheet_name, df in sheets.items():
2783 31                    # Replace NaNs with empty strings for cleaner looking
2784 32                    tables
2785 33                    df_clean = df.fillna("")
2786 34
2787 35                    # We limit the rows to avoid overwhelming the context
2788 36                    window
2789 37                    preview_df = df_clean.head(preview_rows)
2790 38
2791 39                    data_output[sheet_name] = preview_df.to_markdown(index=
2792 40                False)
2793 41            else:
2794 42                return {"error": "Unsupported file format."}
2795 43
2796 44            # Step 4: Return result
2797 45            return {
2798 46                "file_name": os.path.basename(file_path),
2799 47                "sheet_content": data_output,
2800 48                "note": "Visual styles (colors, fonts) were not extracted."
2801 49            }
2802
2803    except Exception as e:
2804        return {"error": f"Failed to parse file: {str(e)}"}

```

Code 5: Code Implenmenation for Read Excel Basic

2800  
2801  
2802  
2803  
2804  
2805  
2806  
2807