



AGENT KB: Leveraging Cross-Domain Experience for Agentic Problem Solving

Xiangru Tang[¶][○], Tianrui Qin[○], Tianhao Peng[○], Ziyang Zhou[¶], Daniel Shao[¶], Tingting Du[¶],
Xinming Wei[○], Peng Xia^{*}, Fang Wu[¶], He Zhu[○], Ge Zhang[▲], Jiaheng Liu[■],
Xingyao Wang[□], Sirui Hong[□], Chenglin Wu[□], Hao Cheng[■], Chi Wang[□], Wangchunshu Zhou[○]

[¶]Yale University, [○]OPPO, [■]UW-Madison, ^{*}UNC Chapel Hill,
[▲]Stanford University, [▲]Bytedance, [■]Nanjing University, [□]All Hands AI,
[□]DeepWisdom, [■]Microsoft Research, [□]Google DeepMind

DOI: <https://github.com/OPPO-PersonalAI/Agent-KB>

Abstract

Current AI agents cannot effectively learn from each other’s problem-solving experiences or use past successes to guide self-reflection and error correction in new tasks. We introduce AGENT KB, a shared knowledge base that captures both high-level problem-solving strategies and detailed execution lessons, enabling knowledge transfer across agent frameworks. Agent KB implements a novel teacher-student dual-phase retrieval mechanism where student agents retrieve workflow-level patterns for strategic guidance while teacher agents identify execution-level patterns for refinement. This hierarchical approach enables agents to break out of limited reasoning pathways by incorporating diverse strategies from external sources. Evaluations on the GAIA benchmark demonstrate substantial performance gains, with AGENT KB improving success rates by up to 6.06 percentage points overall under pass@1. For SWE-bench code repair tasks, our system significantly improved resolution rates, with o3-mini achieving an 8.67 percentage point gain (23 percent to 31.67 percent) in pass@1. Our ablation studies demonstrate that the refinement module proves most critical, with its removal causing 3.85% drop on challenging Level 3 tasks, highlighting that effective knowledge transfer necessitates both strategic guidance and execution-level refinement.

1 Introduction

As artificial intelligence advances, language agents are becoming increasingly vital for solving complex problems [1–6]. Although these agents have demonstrated impressive capabilities through supervised learning and reinforcement learning, they continue to struggle with complex long-term tasks that require sophisticated planning and tool use [7–9]. The integration of autonomous improvement modules has demonstrated performance gains [10–17], yet a critical bottleneck persists.

The fundamental limitation lies in error correction during complex reasoning. When agents encounter difficulties, self-feedback proves insufficient—they lack access to the diverse reasoning strategies and implicit reward signals that guide human experts. Recent work [12, 18–20] shows that *learning reusable experiences (or referred to as memories) from past explorations improves performance*. However, current approaches remain limited to task-specific experiences that operate in isolation. This

*Equal contribution.

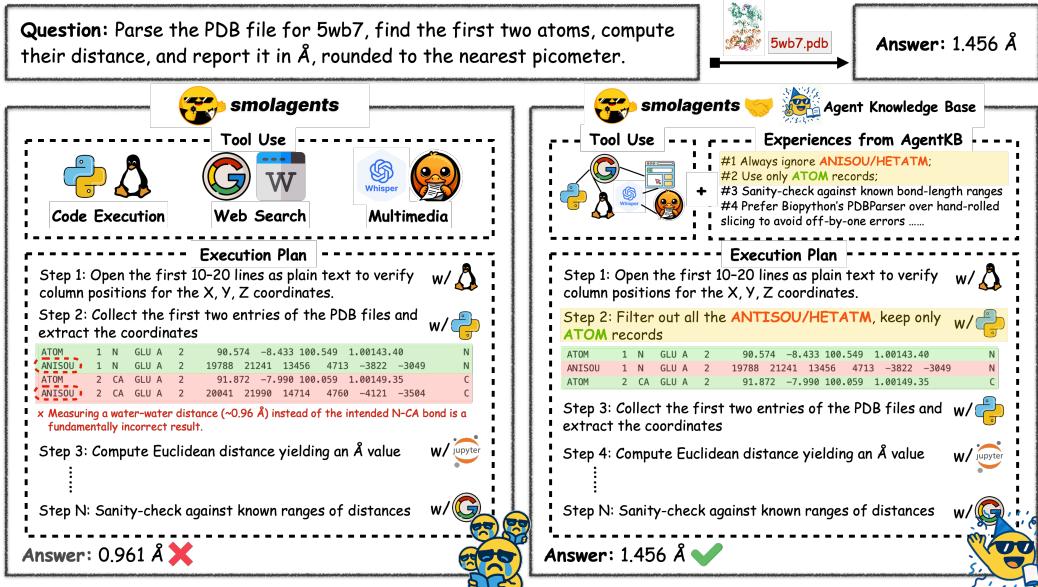


Figure 1: **Comparison of PDB distance-calculation workflows with and without AGENT KB.** **(A) Original pipeline:** indiscriminately reads the first two ATOM/HETATM/ANISOU lines, often selecting solvent records and yielding a spurious O–H distance (0.961 Å). **(B) AGENT KB-enhanced agent workflow:** applies experience-driven rules—filter out all ANISOU/HETATM, use only genuine ATOM entries in file order, and sanity-check against known N–CA bond-length ranges—to correctly extract the backbone N–CA pair and report the distance of 1.456 Å.

isolation forces agents to rediscover similar problem-solving strategies repeatedly when encountering new task types, even when successful approaches from related domains could be adapted and reused.

To understand why current approaches fall short, we identify three critical design flaws in agent experience systems: **(1) Task-Specific Experience Isolation**—agents struggle to transfer knowledge across different task types, forcing them to start from scratch when encountering new domains without leveraging successful strategies from related task categories. **(2) Single-Level Retrieval Granularity**—existing systems employ uniform retrieval mechanisms that fail to distinguish between different phases of problem-solving. Initial planning requires high-level workflow guidance retrieved based on problem characteristics, while mid-execution refinement needs fine-grained reasoning corrections retrieved based on current execution traces. **(3) Static Experience Replay**—systems store and reuse experiences in their original form without proper abstraction, preventing effective adaptation to new contexts where abstract principles would transfer more effectively than raw execution details.

We propose the Agent Knowledge Base (AGENT KB), a hierarchical experience framework that enables complex agentic problem solving through our novel **Reason-Retrieve-Refine** pipeline. Unlike existing systems, AGENT KB enables agents to learn from diverse problem-solving strategies and generalize these experiences across different tasks and frameworks. Specifically, AGENT KB first engages agents in preliminary reasoning about the problem, directing subsequent experience retrieval to relevant solution patterns rather than just matching surface features. We design a novel teacher-student dual-phase retrieval mechanism: student agents first retrieve workflow-level patterns to structure their approach, while teacher agents subsequently identify specific execution patterns to refine implementation details. This hierarchical process enables agents to break out of their limited reasoning pathways by incorporating diverse reasoning strategies from external sources, providing implicit reward signals that guide refinement toward successful solutions.

Our experimental evaluations on the GAIA benchmark demonstrate substantial performance gains. Under pass@1 evaluation, AGENT KB improves success rates by up to 6.06 percentage points overall, with GPT-4.1 showing improvement from 55.15% to 61.21%. Under pass@3 evaluation, the +AGENT KB ✓ (as defined in Section 4.1) configuration achieves even more impressive results, with GPT-4.1 improving from 68.48% to 73.94% and Claude-3.7 from 72.73% to 75.15%. Notably, on challenging Level 3 tasks, Claude-3.7 with +AGENT KB ✓ shows substantial improvement from 50.00% to 57.69%. For SWE-bench code repair tasks, our system significantly improved resolution rates, with o3-mini achieving an 8.67 percentage point gain (23% to 31.67%) under pass@1. Our ablation studies reveal that the hybrid retrieval approach outperforms both pure text similarity and

semantic similarity methods. Notably, automatically generated knowledge achieves comparable performance to manually crafted examples (75.15% vs 76.97% on GAIA), highlighting the value of our knowledge acquisition pipeline.

2 Related Work

2.1 Memory Systems in LLM Agents

Memory systems in LLM agents have evolved from simple storage mechanisms to sophisticated architectures supporting complex reasoning [4, 21–23]. Early implementations like MemoryLLM [24] embedded knowledge in the latent space, while subsequent approaches introduced structured organization through Zettelkasten-style-graph-based systems (A-MEM [15], AriGraph [25]) and hierarchical frameworks (MemGPT [26], Unified Mind Model [27]). Knowledge integration approaches address planning capabilities and hallucination mitigation through frameworks such as Agent Workflow Memory [12], which enables the automatic induction and reuse of sub-workflows, and KnowAgent [28], which augments prompts with action-knowledge bases. More sophisticated approaches include parametric world-knowledge models (WKM) [29] and multi-agent adaptation systems MARK [30]. EcoAssistant [23] demonstrated the effectiveness of knowledge reuse and transfer across agents, establishing a foundation for collaborative reasoning. ReAct [31] synergizes reasoning and acting by interleaving chain-of-thought with tool calls, allowing real-time plan adaptation, while Reflexion [32] enables agents to learn from verbalized self-critiques. Toolformer [33] demonstrates that LLMs can learn to use external tools in an unsupervised manner, patching capability gaps mid-execution. Retrieval mechanisms for memory have progressed beyond basic RAG paradigms [34], with innovations like HippoRAG’s [35] hippocampal-inspired indexing, Echo’s [36] temporal cues, and HiAgent’s [13] sub-goal chunking.

2.2 Multi-Agent Collaboration and Shared Memory

Most existing memory systems remain agent-specific, designed for recalling interaction history [37], modeling user preferences [38], and etc. Memory-augmented embodied agents [39] have begun to explore collaborative architectures, where specialized agents (routing, planning, knowledge base) work together, leveraging in-context learning and RAG to retrieve context from past interactions. However, these systems typically maintain separate memory structures rather than a unified knowledge ecosystem. Limited work exists on cross-agent experience reuse and adaptation. Synapse [10] introduces exemplar memory for trajectory storage but primarily focuses on single-agent contexts. EventWeave [40] addresses incomplete context tracking by identifying both core and supporting events in a dynamic event graph but doesn’t fully extend to multi-agent scenarios. Some researchers have explored pre-conditions for memory-learning agents [14], revealing that memory induction quality significantly impacts performance. This suggests that creating high-quality shared memory structures could benefit multiple agents simultaneously, particularly if stronger agents can induce memories that weaker agents can later leverage. Case-Based Reasoning (CBR) approaches [41] provide promising directions for multi-agent experience reuse, as they enable solving new problems by referencing past experiences.

3 Methodology

As shown in Figure 2, AGENT KB consists of two main phases: AGENT KB *construction* and AGENT KB *enhanced inference*. In the construction phase (left side), we extract generalizable experiences from raw execution logs collected across multiple datasets. During the enhanced inference phase (right side), when facing a new task, an execution agent performs the actual task solving while two additional agents, a student agent and a teacher agent, implement our novel **Reason-Retrieve-Refine** pipeline. These auxiliary agents retrieve relevant experiences from AGENT KB and adaptively refine them to provide targeted guidance that enhances the execution agent’s reasoning and problem-solving capabilities.

3.1 AGENT KB Construction

Before detailing our construction process, we establish the core concepts used throughout our framework. *Experiences* are structured, abstracted problem-solving patterns extracted from raw

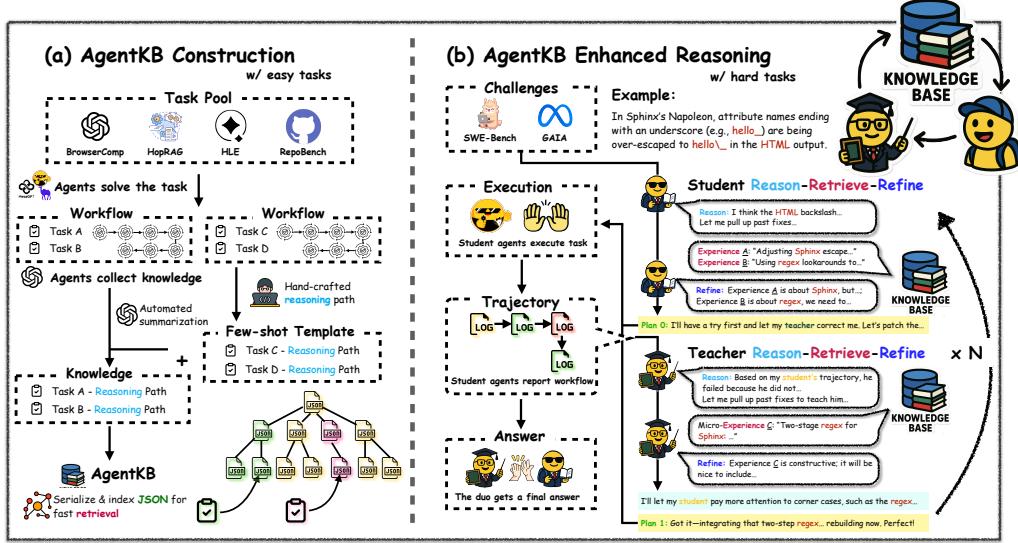


Figure 2: System architecture of AGENT KB, showing the integration of knowledge abstraction, dual-phase retrieval, and adaptive refinement into a unified framework. The student agent retrieves workflow-level patterns for structuring the approach, while the teacher agent retrieves step-level patterns for execution precision.

execution logs, formally represented as tuples $\mathcal{E} = \langle \pi, \gamma, \mathcal{S}, \mathcal{C}, \mathcal{R} \rangle$, containing problem patterns, goals, solution trajectories, context, and relationships. π represents the problem; γ denotes an optional goal or objective; $\mathcal{S} = \{s_1, s_2, \dots, s_n\}$ is an abstracted solution trajectory with reasoning templates, optionally with observed failure modes; \mathcal{C} captures problem characteristics such as domain and difficulty level (optional); \mathcal{R} encodes relational links to other experiences within the hierarchical structure of AGENT KB, though it is excluded in the current implementation to maintain architectural simplicity. These experiences are derived from trajectories—complete execution paths of agents solving tasks, including intermediate steps, tool calls, reasoning traces, and outcomes. At a higher level, workflows capture strategic problem-solving sequences including tool selection, reasoning steps, and decision points.

Our construction process draws from diverse task pools, which are collections of problem-solving scenarios from multiple domains (e.g., BrowseComp, HopRAG, SWE-Bench) that serve as source data for experience extraction. The resulting knowledge forms a collective repository of experiences stored in AGENT KB, organized hierarchically to enable cross-domain transfer and reuse. During inference, the system processes challenges (new tasks)—specific problems or queries that require the retrieval and application of relevant experiences from AGENT KB.

The abstraction process involves both manual inspection of failure cases and automated synthesis using a structured template-based approach to generalize planning behaviors across tasks and domains. Our fundamental hypothesis is that experiences abstracted from simpler tasks can be generalized and adapted to novel, more complex problems by capturing high-level planning patterns rather than framework-specific implementations (see Appendix A).

The process begins by collecting execution logs from previously completed tasks. These logs include both successful and failed trajectories, with special attention paid to common error patterns that hinder performance. To guide generalization, we incorporate manually annotated failure cases and their corresponding resolution strategies as few-shot examples. These exemplars are used to prompt an LLM-based experience generator to produce structured entries in our standardized format, enabling the scalable creation of abstracted experiences across diverse tasks and domains.

Rather than storing raw execution logs, AGENT KB maintains abstracted reasoning patterns that capture generalizable problem-solving strategies, creating a more efficient and transferable knowledge structure.

3.2 Teacher-Student Dual-Phase Inference

We implement a hierarchical teacher-student framework, where both agents operate using complementary **Reason-Retrieve-Refine** (RRR) cycles to solve complex tasks. In this architecture, the execution agent performs the actual task solving, while the student agent retrieves workflow-level experiences to provide initial strategic guidance, and the teacher agent analyzes execution trajectories to identify errors and retrieve step-level experiences for refinement. The teacher supervises the student by detecting reasoning flaws and providing corrective feedback based on relevant experiences from AGENT KB. Notably, both auxiliary agents focus on experience retrieval and adaptive refinement to enhance the performance of the execution agent that carries out the actual problem-solving tasks.

Algorithm 1 Student Agent Inference

```

1: Input: Query  $Q$ , Knowledge base  $\mathcal{K}$ 
2: Output: Execution trajectory  $\mathcal{S}$ , Plan  $\Pi$ 
3: function STUDENTINFERENCE( $Q, \mathcal{K}$ )
4:    $\hat{\pi}, \hat{\gamma} \leftarrow \text{PARSEQUERY}(Q)$                                  $\triangleright$  Extract problem pattern and goal
5:   /* Reason */
6:    $\mathcal{T} \leftarrow \text{REASON}(Q, \hat{\pi}, \hat{\gamma})$                                  $\triangleright$  Initial reasoning
7:   /* Retrieve */
8:    $\mathcal{E}_w \leftarrow \text{top-k}_{\mathcal{E}_i \in \mathcal{K}} [\alpha \cdot \phi_r(\mathcal{E}_i, \mathcal{T}, \hat{\pi}, \hat{\gamma})]$ 
9:   /* Refine */
10:   $\Pi \leftarrow \text{REFINEANDPLAN}(\mathcal{E}_w, \mathcal{T})$                                  $\triangleright$  Integrate and construct plan
11:   $\mathcal{S} \leftarrow \text{EXECUTE}(\Pi)$                                                $\triangleright$  Execute plan
12:  return  $\mathcal{S}, \Pi$ 
13: end function

```

In the **student phase**, the student agent first analyzes query Q to identify the problem ($\hat{\pi}$) and goal ($\hat{\gamma}$), generating initial thoughts \mathcal{T} about potential solutions. Next, it retrieves relevant workflow patterns from AGENT KB:

$$\mathcal{E}_w = \text{top-k}_{\mathcal{E}_i \in \mathcal{K}} [\alpha \cdot \phi_r(\mathcal{E}_i, \mathcal{T}, \hat{\pi}, \hat{\gamma})],$$

where \mathcal{K} is AGENT KB, ϕ_r measures relevance, and α, β are weights. The student then refines these workflows by integrating them with initial reasoning to create and execute a structured plan, resulting in a series of reasoning steps. Knowledge is dynamically adapted rather than directly copied, enabling effective transfer even between dissimilar domains.

Algorithm 2 Teacher Agent Inference

```

1: Input: Query  $Q$ , Execution trajectory  $\mathcal{S}$ , Plan  $\Pi$ , Knowledge base  $\mathcal{K}$ 
2: Output: Guidance  $\Gamma$ 
3: function TEACHERINFERENCE( $Q, \mathcal{S}, \Pi, \mathcal{K}$ )
4:   /* Reason */
5:    $\mathcal{Z} \leftarrow \text{SUMMARIZETRAJECTORY}(\mathcal{S})$                                  $\triangleright$  Summarize trajectory
6:   /* Retrieve */
7:    $\mathcal{E}_s \leftarrow \text{top-m}_{\mathcal{E}_j \in \mathcal{K}} \sum_{s_i \in \mathcal{Z}} [\alpha \cdot \phi_r(s_i, \mathcal{S}_j)]$ 
8:   /* Refine */
9:    $\mathcal{E}_p \leftarrow \text{TRANSFEREXPERIENCE}(\mathcal{E}_s, \mathcal{Z})$                                  $\triangleright$  Transfer relevant experience based on target trajectory
10:   $\Gamma \leftarrow \text{FORMULATEGUIDANCE}(\mathcal{E}_p, Q)$                                  $\triangleright$  Create targeted guidance
11:  return  $\Gamma$ 
12: end function

```

In the **teacher phase**, the teacher agent evaluates the student's reasoning steps by summarizing them and identifying errors along with their types and causes. It retrieves targeted step-level experiences from AGENT KB to address these execution issues:

$$\mathcal{E}_s = \text{top-m}_{\mathcal{E}_j \in \mathcal{K}} \sum_{s_i \in \mathcal{Z}} [\alpha \cdot \phi_r(s_i, \mathcal{S}_j)],$$

where ϕ_r measures similarity. The teacher refines these step-level patterns into precise guidance, providing targeted interventions. This iterative feedback loop progressively enhances the student's performance.

4 Experiment

4.1 Setup

Datasets Our evaluation utilizes two representative benchmarks that assess the diverse capabilities of agents. The GAIA benchmark [42] provides a comprehensive evaluation framework for general AI assistants, containing 165 evaluation instances carefully stratified across three difficulty levels: 53 tasks in Level 1 (basic), 86 tasks in Level 2 (intermediate), and 26 tasks in Level 3 (advanced). These tasks span information retrieval, multi-step reasoning, and complex problem-solving scenarios. The SWE-bench [7] serves as our second benchmark, focusing on realistic software engineering challenges extracted from GitHub issues, requiring agents to understand existing codebases and implement appropriate fixes.

The knowledge base for AGENT KB draws from diverse sources. For general assistant tasks, we aggregate experiences from four complementary datasets: BrowseComp [43] (1,266 tasks), HopRAG [44] (2,556 tasks), a text-based subset of HLE [45] (3,000 tasks), and WebWalkerQA [46] (680 tasks). For software engineering knowledge, we incorporate structured experiences from three major repositories: RepoClassBench [47], SWE-Gym-Raw [48], and RepoEval [49], collectively comprising approximately 3,000 structured problem-solving traces (see Appendix I).

Model Configurations We evaluate three distinct configurations across multiple foundation models to assess the effectiveness of AGENT KB. We use an improved version of the smolagents² [50] to serve as our base agent framework. The framework we use extends the smolagents with audio-visual comprehension modules and a multi-source retrieval system, enhancing multimodal input processing and enabling more efficient access to heterogeneous information sources. For the SWE-bench benchmark, we employ the OpenHands framework³ as our base agent framework. Default settings are used for all hyperparameters unless otherwise noted. The +AGENT KB configuration implements a two-round, teacher-student knowledge adaptation process: first, the execution agent attempts to solve the task; then, the student and teacher agents find relevant experiences and provide feedback without knowing whether the execution agent’s solution was correct (unsupervised). The execution agent then makes a second attempt to incorporate this feedback. In our +AGENT KB ✓ configuration, we will only improve the examples that fail in the first round of execution. Thus, we conduct a comparison between this setup and the smolagents baseline under the pass@2 evaluation. Same as before, the teacher and student agents first analyze the initial execution trace, which consists of reasoning steps, tool calls, and intermediate results, to detect potential flaws. Upon identifying such errors, the agents retrieve relevant prior experiences and refine them through contextual alignment and error-specific adaptation. Since the smolagents framework is a simple one, to test an upper bound of our approach and ensure a comparison with current complex methods that employ various performance-enhancing tricks, we implement similar improvements in +AGENT KB ✓✓. It includes optimized retrieval mechanisms, fine-grained knowledge extraction patterns, consistent output formatting corrections, and utilizes the pass@3 score. The specific implementation details of the +AGENT KB ✓✓ configuration and computational cost can be found in the Appendix J.2.

The evaluated LLMs include GPT-4o(2024-11-20), GPT-4.1(gpt-4.1-2025-04-14), Claude-3.7(claude-3-7-sonnet-20250219), o3-mini(o3-mini-2025-01-31), Qwen-3 32B, and DeepSeek-R1. In all experiments, we maintained consistent hyperparameters across comparable settings, with top_k set to 0.1 and temperature set to 1.0.

4.2 Main Results

In Table 1, our approach demonstrates significant improvements over baselines across all GAIA’s difficulty levels. GPT-4.1 with +AGENT KB ✓✓ shows an overall improvement of 18.79 percentage points, with the largest gains (19.77 points) observed in medium-difficulty tasks (Level 2). Claude models exhibit similar benefits from AGENT KB integration, with Claude-3.7 with +AGENT KB ✓✓ improving from 58.79% to 75.15% in overall performance. Figure 3 also demonstrates consistent performance improvements across all six base LLMs tested. A 19.23 percentage point gain (Claude-3.7 rising from 38.46% to 57.69%) in the most complex scenario category (level 3) validates our approach’s effectiveness in supporting sophisticated multi-step reasoning and planning. Such improvements

²<https://github.com/huggingface/smolagents>

³<https://github.com/All-Hands-AI/OpenHands>

Table 1: Performance of various agent frameworks on GAIA benchmark (validation set).

Method	Models	Average	Level 1	Level 2	Level 3
<i>Single Model</i>					
Search-o1-32B [51]	-	39.8	53.8	34.6	16.7
WebThinker-32B-RL [52]	-	48.5	56.4	50.0	16.7
<i>Closed-source Agent Frameworks</i>					
TraseAgent [53]	Claude <i>etc.</i>	70.30	83.02	69.77	46.15
OpenAI Deep Research [54]	Unknown	67.36	74.29	69.06	47.60
h2oGPTe [55]	Claude-3.5	63.64	67.92	67.44	42.31
Desearch [56]	GPT-4o	56.97	71.70	58.14	23.08
Alita [57] (pass@1)	Claude 3.7+GPT-4o	72.73	81.13	75.58	46.15
Alita [57] (pass@2)	Claude 3.7+GPT-4o	78.79	88.68	80.23	53.85
Alita [57] (pass@3)	Claude 3.7+GPT-4o	86.06	96.23	86.04	65.38
<i>Open-Source Agent Frameworks</i>					
OWL Workforce (pass@3) [58]	GPT-4o+o3-mini	60.61	81.14	58.14	26.92
OWL Role Playing (pass@3) [58]	GPT-4o+o3-mini	58.18	81.14	54.65	23.08
TapeAgents [59]	Claude 3.7 <i>etc.</i>	55.76	71.70	53.49	30.77
AutoAgent [60]	Claude 3.5 <i>etc.</i>	55.15	71.70	53.40	26.92
smolagents [61]	GPT-4.1	55.15	67.92	53.49	34.62
Magnetic-1 [62]	OpenAI o1 <i>etc.</i>	46.06	56.60	46.51	23.08
FRIDAY [63]	GPT-4 turbo	34.55	45.28	34.88	11.54
<i>Pass@1</i>					
smolagents (baseline)	GPT-4.1	55.15	67.92	53.49	34.62
smolagents +AGENT KB	GPT-4.1	61.21 <small>±6.06</small>	79.25 <small>±11.33</small>	58.14 <small>±4.65</small>	34.62
smolagents (baseline)	Claude 3.7 <i>etc.</i>	58.79	64.15	61.63	38.46
smolagents +AGENT KB	Claude 3.7 <i>etc.</i>	65.45 <small>±6.66</small>	75.47 <small>±11.32</small>	66.28 <small>±4.65</small>	38.46
<i>Pass@2</i>					
smolagents (baseline)	GPT-4.1	61.82	73.58	62.79	34.62
smolagents +AGENT KB ✓	GPT-4.1	67.27 <small>±5.45</small>	83.02 <small>±9.44</small>	67.44 <small>±4.65</small>	34.62
smolagents (baseline)	Claude 3.7 <i>etc.</i>	63.64	77.36	61.63	42.31
smolagents +AGENT KB ✓	Claude 3.7 <i>etc.</i>	69.70 <small>±6.06</small>	79.25 <small>±11.89</small>	69.77 <small>±8.14</small>	50.00 <small>±7.69</small>
<i>Pass@3</i>					
smolagents (baseline)	GPT-4.1	68.48	77.36	68.60	50.00
smolagents +AGENT KB ✓♡	GPT-4.1	73.94 <small>±5.46</small>	84.91 <small>±7.55</small>	73.26 <small>±4.66</small>	53.85 <small>±3.85</small>
smolagents (baseline)	Claude 3.7 <i>etc.</i>	72.73	81.13	74.42	50.00
smolagents +AGENT KB ✓♡	Claude 3.7 <i>etc.</i>	75.15 <small>±2.42</small>	84.91 <small>±3.78</small>	74.42	57.69 <small>±7.69</small>

suggest that the bottleneck in handling complex tasks lies in their ability to effectively leverage relevant past experiences. More results can be found in Appendix 9.

Notably, the +AGENT KB ✓♡ -enhanced Claude-3.7 model achieves an average GAIA score of 75.15%, surpassing closed-source systems like h2oGPTe (63.64%) and open-source frameworks like OWL (69.09%). This performance is particularly impressive given that our approach builds upon a relatively straightforward agent framework (smolagents).

For the SWE-bench lite benchmark [7], we set the max limit for agent iterations to 50 and 100, and conduct experiments, respectively. Table 2 shows similar patterns of improvement across different model types. Claude-3.7 achieves the most substantial gains, with performance increasing from 30.00% to 51.00% at 50 iterations. Interestingly, we observe that the relative magnitude of improvement correlates with model sophistication, with larger and more capable models, such as GPT-4.1, showing more substantial gains than smaller models, like Qwen-3 32B. This suggests that more advanced models are better able to leverage the retrieved knowledge, potentially due to their enhanced reasoning capabilities.

4.3 Ablation Studies

To assess the contribution of each core component in AGENT KB, we conduct systematic ablation studies (Table 3). All experiments are performed on +AGENT KB . For detailed ablation configurations and experimental setup, please refer to Appendix D.3.

Removing either the student or teacher agent reduces performance to 59.39%, highlighting the complementary roles of both in the dual-phase architecture. Notably, the student agent is especially important for Level 1 tasks (a drop from 79.25% → 75.47%), suggesting its key role in planning simpler workflows. In contrast, removing the teacher agent leads to a sharper decline in Level 1

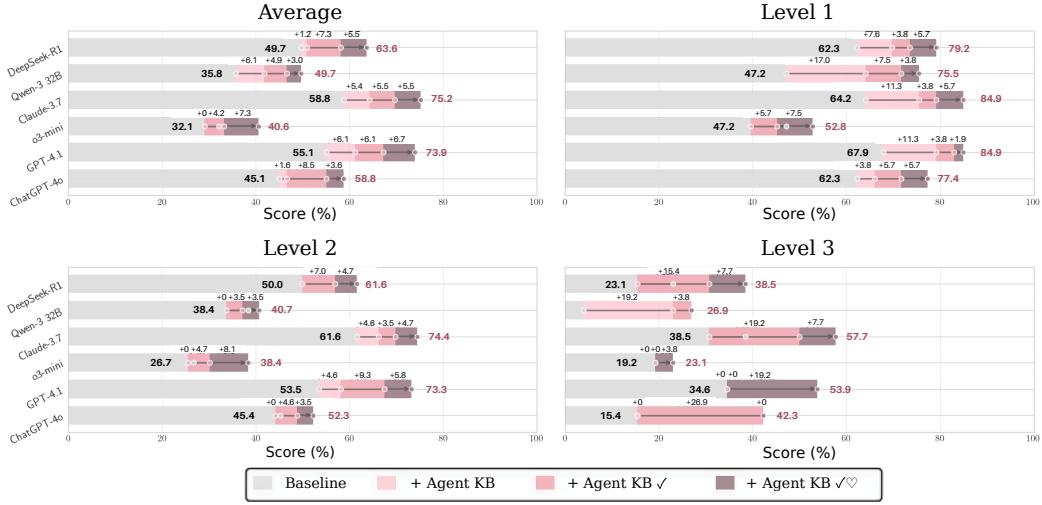


Figure 3: Score improvements (%) across difficulty levels for multiple base LLMs enhanced with AGENT KB.

Table 2: Main results on the SWE-bench lite with maximum iteration limits of 50 and 100.

Method	Models	Success Rate	
		Max Iter 50	Max Iter 100
OpenHands (baseline)		16.33	26.00
OpenHands +AGENT KB (pass@1)	GPT-4o	20.33 ↑+4.00	29.67 ↑+3.67
OpenHands +AGENT KB ✓ (pass@2)		29.33	35.67
OpenHands +AGENT KB ✓✓ (pass@3)		31.33	39.33
OpenHands (baseline)		24.33	28.67
OpenHands +AGENT KB (pass@1)	GPT-4.1	28.33 ↑+4.00	31.67 ↑+3.00
OpenHands +AGENT KB ✓ (pass@2)		37.33	42.33
OpenHands +AGENT KB ✓✓ (pass@3)		38.67	45.67
OpenHands (baseline)		23.00	29.33
OpenHands +AGENT KB (pass@1)	o3-mini	31.67 ↑+8.67	33.67 ↑+4.34
OpenHands +AGENT KB ✓ (pass@2)		35.33	36.33
OpenHands +AGENT KB ✓✓ (pass@3)		37.00	40.00
OpenHands (baseline)		24.33	30.00
OpenHands +AGENT KB (pass@1)	DeepSeek-R1	26.67 ↑+2.34	33.33 ↑+3.33
OpenHands +AGENT KB ✓ (pass@2)		31.00	35.67
OpenHands +AGENT KB ✓✓ (pass@3)		32.67	37.33
OpenHands (baseline)		18.33	25.67
OpenHands +AGENT KB (pass@1)	Qwen-3 32B	20.67 ↑+2.34	28.67 ↑+3.00
OpenHands +AGENT KB ✓ (pass@2)		28.67	34.33
OpenHands +AGENT KB ✓✓ (pass@3)		30.33	36.67

accuracy (79.25% → 73.58%), indicating its role in early-stage refinement. The most significant drop occurs when the `Refine` module is removed, decreasing overall accuracy by 6.06 percentage points (61.21% → 55.15%) and Level 3 performance by 3.85 points (34.62% → 30.77%), underscoring the necessity of fine-grained error correction. Ablating the `Retrieve` module also yields notable degradation (-3.63 points), demonstrating that knowledge grounding via retrieval is essential. In contrast, omitting the `Reason` module results in only a modest drop (-1.21), suggesting that retrieval and refinement can partially compensate for the absence of high-level planning. Finally, replacing structured experiences with raw workflow logs reduces performance to 58.18%, reaffirming the importance of knowledge abstraction and reuse beyond naive trajectory replay. These results validate that reasoning, retrieval, and refinement each contribute distinct and synergistic improvements, with the refinement phase playing a particularly critical role in ensuring the correctness of execution on challenging tasks.

To better understand the factors contributing to AGENT KB’s effectiveness, we conduct an in-depth analysis of different retrieval strategies across abstraction levels (Figure 4). Using GPT-4.1 as our base

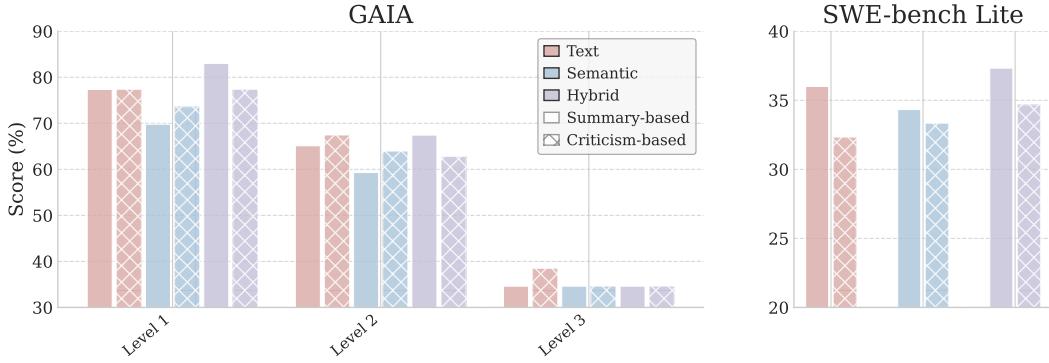


Figure 4: Performance comparison of text, semantic, and hybrid retrieval methods across two different abstraction levels. The left panels display results for summary-based retrieval, while the right panels show results for criticism-based retrieval.

model with top-k=3, we compare three retrieval approaches (text similarity, semantic similarity, and hybrid retrieval) across two complementary abstraction methods that we integrate into our full system. Our implemented system combines both summary-based and criticism-based approaches. Following the RRR pipeline, we first retrieve relevant experiences based on a summary of execution logs generated through reasoning and analysis. We then refine the experiences using two complementary strategies: the summary-based method, which distills logs into concise representations, and the criticism-based approach, which leverages teacher agents to identify and correct potential errors. The refined experiences are subsequently used for the next run. Figure 4 demonstrates their distinct contributions.

For summary-based retrieval (left panels), hybrid methods consistently outperform single-approach strategies, achieving 83% accuracy on Level 1 GAIA tasks and 37% on SWE-bench lite. The performance advantage is particularly pronounced for Level 1 and Level 2 tasks, where hybrid retrieval shows improvements of up to 9 percentage points over semantic-only approaches. Criticism-based retrieval (right panels) exhibits a different pattern, with text similarity performing competitively for Level 2 tasks (67%) and semantic similarity showing stronger results on the SWE-bench (33%). Hybrid approaches maintain their edge in most scenarios, though with narrower margins.

Table 4 reveals that AGENT KB achieves comparable performance to HAND CRAFTED annotations (created by computer science students) across both benchmarks. While manually annotated experiences perform slightly better on AGENT KB tasks (79.07% vs 74.42%), AGENT KB outperforms human annotations on challenging Level 3 tasks (57.69% vs 53.85%).

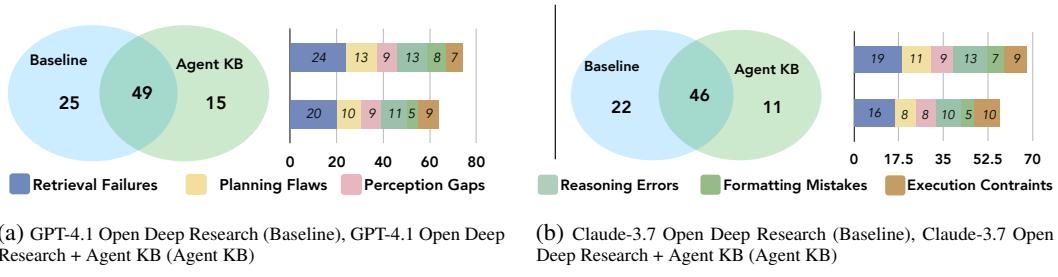
While keeping the agent framework architecture (GPT-4.1) identical, we vary the underlying LLMs for both teacher and student agents.

Table 3: Ablation results for different components of the AGENT KB.

Ablation Setting	Avg	Level 1	Level 2	Level 3
smolagent	55.15	67.92	53.49	34.62
smolagents +AGENT KB	61.21	79.25	58.14	34.62
w/o Student Agent	59.39	75.47	56.98	34.62
w/o Teacher Agent	59.39	73.58	58.14	34.62
w/o Reason Module	60.00	77.36	56.98	34.62
w/o Retrieve Module	57.58	73.58	54.65	34.62
w/o Refine Module	55.15	69.81	53.49	30.77
w/ Raw Workflow	58.18	73.58	55.81	34.62

Table 4: Performance comparison of different knowledge types across GAIA and SWE-bench benchmarks. Baseline shows results without experience augmentation, while HAND CRAFTED represents experiences manually annotated by three computer science students. AGENT KB shows results with automatically extracted and refined experiences.

Knowledge type	Average	GAIA			SWE-bench Lite
		Level 1	Level 2	Level 3	
Baseline	55.15	67.92	53.49	34.62	24.33
+ HAND CRAFTED	76.97	84.91	79.07	53.85	55.67
+ AGENT KB	75.15	84.91	74.42	57.69	51.00



(a) GPT-4.1 Open Deep Research (Baseline), GPT-4.1 Open Deep Research + Agent KB (Agent KB) (b) Claude-3.7 Open Deep Research (Baseline), Claude-3.7 Open Deep Research + Agent KB (Agent KB)

Figure 5: The frequency of errors with and without **Agent KB**. The Venn diagrams quantify overlapping and unique failure cases, while the horizontal bar charts show category-specific error counts.

In Table 5, the results show that Claude-3.7-sonnet achieves slightly higher performance on both GAIA (~70%) and SWE-bench lite (~39%), followed by GPT-4.1, GPT-4o, and o3-mini. To facilitate a fair comparison of the capabilities of different teacher agents based on their backbone models, we standardize the evaluation using AGENT KB ✓ throughout this analysis. This allows us to isolate the impact of the underlying LLM while keeping the agent framework constant. While more sophisticated models provide marginal benefits, AGENT KB’s performance improvements remain robust across different underlying LLMs. Notably, even when using GPT-4.1 for both teacher and student agents (matching our baseline’s underlying model), we observe substantial performance gains. This demonstrates that AGENT KB enables a form of self-correction where the same LLM can effectively critique and improve its work through the teacher-student architecture.

4.4 Error Analysis

For GPT-4.1 (Figure 5a), we observe that 49 errors occur in both configurations, while 25 errors specific to the baseline were successfully corrected by AGENT KB. The enhanced model introduced only 15 new errors, yielding a net error reduction of 10 instances. Similarly, with Claude-3.7 (Figure 5b), 46 errors persist across both configurations, while AGENT KB corrects 22 baseline-specific errors and introduces just 11 new ones, resulting in a net improvement of 11 instances. The bar charts reveal the distribution of error types. The authors manually reviewed and categorized each error case through a systematic annotation process to ensure accurate classification across six distinct error categories. For GPT-4.1, retrieval errors decreased from 24 to 20 instances, and planning errors from 13 to 10. Claude-3.7 demonstrates even more pronounced improvements in retrieval (19 to 16) and reasoning errors (13 to 8). This improvement stems from AGENT KB’s knowledge base, which contains analogous search protocols and workflows, allowing agents to accumulate expertise through standardized pathways and successful planning precedents. Formatting errors also decreased significantly as agents adopt format requirements derived from similar experiences, contributing to more precise output specifications. While image and video comprehension tasks remain constrained by underlying tool capabilities, AGENT KB-enhanced agents still formulate more appropriate plans for visual tool utilization. Furthermore, the knowledge base helps reduce task hallucinations, resulting in more streamlined planning steps that minimize context length and information loss during complex reasoning processes. Interestingly, while both models exhibit similar patterns of improvement, Claude-3.7 achieves greater error reduction in reasoning tasks, whereas

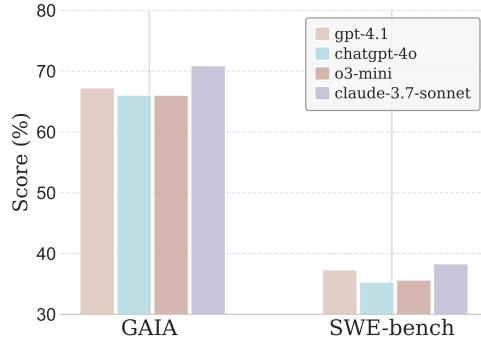


Table 5: Performance comparison of four teacher and student agent models (GPT-4.1, GPT-4o, o3-mini, and Claude-3.7-sonnet) across GAIA and SWE-bench lite benchmarks.

GPT-4.1 benefits more from perception gap resolution, highlighting how AGENT KB’s effectiveness complements each model’s inherent strengths and weaknesses.

Figure 5 illustrates the impact of AGENT KB on error patterns across different base LLM configurations. The Venn diagrams provide a quantitative comparison of errors between the smolagents framework and its AGENT KB-enhanced counterparts.

5 Conclusion

We introduce AGENT KB, a unified and scalable framework that enables LLM agents to continuously learn from experience across tasks, domains, and agent architectures. By structuring prior workflows into generalizable experience units and supporting their reuse through a dual-phase, teacher-student retrieval and refinement pipeline, AGENT KB moves beyond simple memory replay to realize adaptive, experience-driven reasoning. Our experiments across diverse settings—including GAIA and SWE-bench—demonstrate consistent performance improvements across difficulty levels, model families, and agent frameworks. Notably, AGENT KB’s structured knowledge abstraction and dual-phase inference enable not only effective reuse of past solutions but also the evolution of better workflows through agent collaboration. These results position AGENT KB as a general-purpose infrastructure for scalable, continual improvement in agent ecosystems, bridging the gap between episodic memory and cumulative agent intelligence.

References

- [1] Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. Chateval: Towards better llm-based evaluators through multi-agent debate. *arXiv preprint arXiv:2308.07201*, 2023.
- [2] Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, et al. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*, 3(4):6, 2023.
- [3] Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xiangliang Zhang. Large language model based multi-agents: A survey of progress and challenges. *arXiv preprint arXiv:2402.01680*, 2024.
- [4] Bang Liu, Xinfeng Li, Jiayi Zhang, Jinlin Wang, Tanjin He, Sirui Hong, Hongzhang Liu, Shaokun Zhang, Kaitao Song, Kunlun Zhu, et al. Advances and challenges in foundation agents: From brain-inspired intelligence to evolutionary, collaborative, and safe systems. *arXiv preprint arXiv:2504.01990*, 2025.
- [5] Wangchunshu Zhou, Yuchen Eleanor Jiang, Long Li, Jialong Wu, Tiannan Wang, Shi Qiu, Jintian Zhang, Jing Chen, Ruipu Wu, Shuai Wang, et al. Agents: An open-source framework for autonomous language agents. *arXiv preprint arXiv:2309.07870*, 2023.
- [6] Wangchunshu Zhou, Yixin Ou, Shengwei Ding, Long Li, Jialong Wu, Tiannan Wang, Jiamin Chen, Shuai Wang, Xiaohua Xu, Ningyu Zhang, et al. Symbolic learning enables self-evolving agents. *arXiv preprint arXiv:2406.18532*, 2024.
- [7] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*, 2023.
- [8] Xu Huang, Weiwen Liu, Xiaolong Chen, Xingmei Wang, Hao Wang, Defu Lian, Yasheng Wang, Ruiming Tang, and Enhong Chen. Understanding the planning of llm agents: A survey. *arXiv preprint arXiv:2402.02716*, 2024.
- [9] Weimin Xiong, Yifan Song, Qingxiu Dong, Bingchan Zhao, Feifan Song, Xun Wang, and Sujian Li. Mpo: Boosting llm agents with meta plan optimization. *arXiv preprint arXiv:2503.02682*, 2025.
- [10] Longtao Zheng, Rundong Wang, Xinrun Wang, and Bo An. Synapse: Trajectory-as-exemplar prompting with memory for computer control. *arXiv preprint arXiv:2306.07863*, 2023.

- [11] Zeyu Zhang, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Quanyu Dai, Jieming Zhu, Zhenhua Dong, and Ji-Rong Wen. A survey on the memory mechanism of large language model based agents. *arXiv preprint arXiv:2404.13501*, 2024.
- [12] Zora Zhiruo Wang, Jiayuan Mao, Daniel Fried, and Graham Neubig. Agent workflow memory. *arXiv preprint arXiv:2409.07429*, 2024.
- [13] Mengkang Hu, Tianxing Chen, Qiguang Chen, Yao Mu, Wenqi Shao, and Ping Luo. Hiagent: Hierarchical working memory management for solving long-horizon agent tasks with large language model. *arXiv preprint arXiv:2408.09559*, 2024.
- [14] Vishwa Shah, Vishruth Veerendranath, Graham Neubig, Daniel Fried, and Zora Zhiruo Wang. Exploring the pre-conditions for memory-learning agents. In *Scaling Self-Improving Foundation Models without Human Supervision*, 2025.
- [15] Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao, Juntao Tan, and Yongfeng Zhang. A-mem: Agentic memory for llm agents. *arXiv preprint arXiv:2502.12110*, 2025.
- [16] Siwei Liu, Jinyuan Fang, Han Zhou, Yingxu Wang, and Zaiqiao Meng. Sew: Self-evolving agentic workflows for automated code generation. *arXiv preprint arXiv:2505.18646*, 2025.
- [17] Jenny Zhang, Shengran Hu, Cong Lu, Robert Lange, and Jeff Clune. Darwin godel machine: Open-ended evolution of self-improving agents. *arXiv preprint arXiv:2505.22954*, 2025.
- [18] Xiaoyu Tan, Bin Li, Xihe Qiu, Chao Qu, Wei Chu, Yinghui Xu, and Yuan Qi. Meta-agent-workflow: Streamlining tool usage in llms through workflow construction, retrieval, and refinement. In *Companion Proceedings of the ACM on Web Conference 2025*, pages 458–467, 2025.
- [19] Prateek Chhikara, Dev Khant, Saket Aryan, Taranjeet Singh, and Deshraj Yadav. Mem0: Building production-ready ai agents with scalable long-term memory. *arXiv preprint arXiv:2504.19413*, 2025.
- [20] Zijian Zhou, Ao Qu, Zhaoxuan Wu, Sunghwan Kim, Alok Prakash, Daniela Rus, Jinhua Zhao, Bryan Kian Hsiang Low, and Paul Pu Liang. Mem1: Learning to synergize memory and reasoning for efficient long-horizon agents. *arXiv preprint arXiv:2506.15841*, 2025.
- [21] Jinghua Piao, Yuwei Yan, Jun Zhang, Nian Li, Junbo Yan, Xiaochong Lan, Zhihong Lu, Zhiheng Zheng, Jing Yi Wang, Di Zhou, et al. Agentsociety: Large-scale simulation of llm-driven generative agents advances understanding of human behaviors and society. *arXiv preprint arXiv:2502.08691*, 2025.
- [22] Ruihong Zeng, Jinyuan Fang, Siwei Liu, and Zaiqiao Meng. On the structural memory of llm agents. *arXiv preprint arXiv:2412.15266*, 2024.
- [23] Jieyu Zhang, Ranjay Krishna, Ahmed Hassan Awadallah, and Chi Wang. Ecoassistant: Using llm assistants more affordably and accurately. In *ICLR 2024 Workshop on Large Language Model (LLM) Agents*, 2025.
- [24] Yu Wang, Yifan Gao, Xiusi Chen, Haoming Jiang, Shiyang Li, Jingfeng Yang, Qingyu Yin, Zheng Li, Xian Li, Bing Yin, Jingbo Shang, and Julian McAuley. MEMORYLLM: Towards self-updatable large language models. *arXiv preprint arXiv:2402.04624*, 2024.
- [25] Petr Anokhin, Nikita Semenov, Artyom Sorokin, Dmitry Evseev, Mikhail Burtsev, and Evgeny Burnaev. Arigraph: Learning knowledge graph world models with episodic memory for llm agents. *arXiv preprint arXiv:2407.04363*, 2024.
- [26] Charles Packer, Sarah Wooders, Kevin Lin, Vivian Fang, Shishir G. Patil, Ion Stoica, and Joseph E. Gonzalez. Memgpt: Towards LLMs as operating systems. *arXiv preprint arXiv:2310.08560*, 2023.
- [27] Pengbo Hu and Xiang Ying. Unified mind model: Reimagining autonomous agents in the llm era. *arXiv preprint arXiv:2503.03459*, 2025.

- [28] Yuqi Zhu, Shuofei Qiao, Yixin Ou, Shumin Deng, Ningyu Zhang, Shiwei Lyu, Yue Shen, Lei Liang, Jinjie Gu, and Huajun Chen. Knowagent: Knowledge-augmented planning for LLM-based agents. In *Proc. NAACL Findings*, 2024.
- [29] Shuofei Qiao, Runnan Fang, Ningyu Zhang, Yuqi Zhu, Xiang Chen, Shumin Deng, Yong Jiang, Pengjun Xie, Fei Huang, and Huajun Chen. Agent planning with world knowledge model. In *Proc. NeurIPS*, 2024.
- [30] Anish Ganguli, Prabal Deb, and Debleena Banerjee. Mark: Memory augmented refinement of knowledge. *arXiv preprint arXiv:2505.05177*, 2025.
- [31] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. ReAct: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.
- [32] Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *arXiv preprint arXiv:2303.11366*, 2023.
- [33] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*, 2023.
- [34] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Proc. NeurIPS*, 2020.
- [35] Bernal Jiménez Gutiérrez, Yiheng Shu, Yu Gu, Michihiro Yasunaga, and Yu Su. Hipporag: Neurobiologically inspired long-term memory for large language models. *arXiv preprint arXiv:2405.14831*, 2024.
- [36] Baichuan Liu, Chen Li, Ming Tan, Wenqiang Liu, and Yiming Yang. Echo: A large language model with temporal episodic memory. *arXiv preprint arXiv:2502.16090*, 2025.
- [37] Junru Lu, Siyu An, Mingbao Lin, Gabriele Pergola, Yulan He, Di Yin, Xing Sun, and Yunsheng Wu. Memochat: Tuning llms to use memos for consistent long-range open-domain conversation. *arXiv preprint arXiv:2308.08239*, 2023.
- [38] Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. Memorybank: Enhancing large language models with long-term memory. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2024.
- [39] Marc Glockner, Peter Hönig, Matthias Hirschmanner, and Markus Vincze. Llm-empowered embodied agent for memory-augmented task planning in household robotics. *arXiv preprint arXiv:2504.21716*, 2025.
- [40] Zhengyi Zhao, Shubo Zhang, Yiming Du, Bin Liang, Baojun Wang, Zhongyang Li, Binyang Li, and Kam-Fai Wong. Eventweave: A dynamic framework for capturing core and supporting events in dialogue systems. *arXiv preprint arXiv:2503.23078*, 2025.
- [41] Kostas Hatalis, Despina Christou, and Vyshnavi Kondapalli. Review of case-based reasoning for llm agents: Theoretical foundations, architectural components, and cognitive integration. *arXiv preprint arXiv:2504.06943*, 2025.
- [42] Grégoire Mialon, Clémentine Fourrier, Thomas Wolf, Yann LeCun, and Thomas Scialom. Gaia: a benchmark for general ai assistants. In *The Twelfth International Conference on Learning Representations*, 2023.
- [43] Jason Wei, Zhiqing Sun, Spencer Papay, Scott McKinney, Jeffrey Han, Isa Fulford, Hyung Won Chung, Alex Tachard Passos, William Fedus, and Amelia Glaese. Browsecomp: A simple yet challenging benchmark for browsing agents. *arXiv preprint arXiv:2504.12516*, 2025.

- [44] Hao Liu, Zhengren Wang, Xi Chen, Zhiyu Li, Feiyu Xiong, Qinhan Yu, and Wentao Zhang. Hoprag: Multi-hop reasoning for logic-aware retrieval-augmented generation. *arXiv preprint arXiv:2502.12442*, 2025.
- [45] Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang, Mohamed Shaaban, John Ling, Sean Shi, et al. Humanity’s last exam. *arXiv preprint arXiv:2501.14249*, 2025.
- [46] Jialong Wu, Wenbiao Yin, Yong Jiang, Zhenglin Wang, Zekun Xi, Runnan Fang, Linhai Zhang, Yulan He, Deyu Zhou, Pengjun Xie, et al. Webwalker: Benchmarking llms in web traversal. *arXiv preprint arXiv:2501.07572*, 2025.
- [47] Ajinkya Deshpande, Anmol Agarwal, Shashank Shet, Arun Iyer, Aditya Kanade, Ramakrishna Bairi, and Suresh Parthasarathy. Class-level code generation from natural language using iterative, tool-enhanced reasoning over repository. *arXiv preprint arXiv:2405.01573*, 2024.
- [48] Jiayi Pan, Xingyao Wang, Graham Neubig, Navdeep Jaitly, Heng Ji, Alane Suhr, and Yizhe Zhang. Training software engineering agents and verifiers with swe-gym. *arXiv preprint arXiv:2412.21139*, 2024.
- [49] Fengji Zhang, Bei Chen, Yue Zhang, Jacky Keung, Jin Liu, Daoguang Zan, Yi Mao, Jian-Guang Lou, and Weizhu Chen. Repocoder: Repository-level code completion through iterative retrieval and generation. *arXiv preprint arXiv:2303.12570*, 2023.
- [50] He Zhu, Tianrui Qin, King Zhu, Heyuan Huang, Yeyi Guan, Jinxiang Xia, Yi Yao, Hanhao Li, Ningning Wang, Pai Liu, et al. Oagents: An empirical study of building effective agents. *arXiv preprint arXiv:2506.15741*, 2025.
- [51] Xiaoxi Li, Guanting Dong, Jiajie Jin, Yuyao Zhang, Yujia Zhou, Yutao Zhu, Peitian Zhang, and Zhicheng Dou. Search-o1: Agentic search-enhanced large reasoning models. *arXiv preprint arXiv:2501.05366*, 2025.
- [52] Xiaoxi Li, Jiajie Jin, Guanting Dong, Hongjin Qian, Yutao Zhu, Yongkang Wu, Ji-Rong Wen, and Zhicheng Dou. Webthinker: Empowering large reasoning models with deep research capability. *arXiv preprint arXiv:2504.21776*, 2025.
- [53] Trase. Meet trase systems. [Online], 2024. <https://www.trasesystems.com/>.
- [54] OpenAI. deepresearch, 2024.
- [55] H2O.ai. Autonomous agentic ai: execute multi-step workflows autonomously. [Online], 2024. <https://h2o.ai/platform/enterprise-h2ogpte/#AgenticAI>.
- [56] Dsearch AI. Dsearch, 2024.
- [57] Jiahao Qiu, Xuan Qi, Tongcheng Zhang, Xinzhe Juan, Jiacheng Guo, Yifu Lu, Yimin Wang, Zixin Yao, Qihan Ren, Xun Jiang, et al. Alita: Generalist agent enabling scalable agentic reasoning with minimal predefinition and maximal self-evolution. *arXiv preprint arXiv:2505.20286*, 2025.
- [58] Mengkang Hu, Yuhang Zhou, Wendong Fan, Yuzhou Nie, Bowei Xia, Tao Sun, Ziyu Ye, Zhaoxuan Jin, Yingru Li, Zeyu Zhang, Yifeng Wang, Qianshuo Ye, Ping Luo, and Guohao Li. Owl: Optimized workforce learning for general multi-agent assistance in real-world task automation, 2025.
- [59] Dzmitry Bahdanau, Nicolas Gontier, Gabriel Huang, Ehsan Kamalloo, Rafael Pardinas, Alex Piché, Torsten Scholak, Oleh Shliazhko, Jordan Prince Tremblay, Karam Ghanem, Soham Parikh, Mitul Tiwari, and Quaizar Vohra. Tapeagents: a holistic framework for agent development and optimization, 2024.
- [60] Jiabin Tang, Tianyu Fan, and Chao Huang. Autoagent: A fully-automated and zero-code framework for llm agents. *arXiv e-prints*, pages arXiv–2502, 2025.
- [61] LangChain. Open deep research. [Online], 2024. https://github.com/langchain-ai/open_deep_research.

- [62] Adam Fourney, Gagan Bansal, Hussein Mozannar, Cheng Tan, Eduardo Salinas, Friederike Niedtner, Grace Proebsting, Griffin Bassman, Jack Gerrits, Jacob Alber, et al. Magentic-one: A generalist multi-agent system for solving complex tasks. *arXiv preprint arXiv:2411.04468*, 2024.
- [63] Zhiyong Wu, Chengcheng Han, Zichen Ding, Zhenmin Weng, Zhoumianze Liu, Shunyu Yao, Tao Yu, and Lingpeng Kong. Os-copilot: Towards generalist computer agents with self-improvement. *arXiv preprint arXiv:2402.07456*, 2024.
- [64] Yiming Du, Wenyu Huang, Danna Zheng, Zhaowei Wang, Sebastien Montella, Mirella Lapata, Kam-Fai Wong, and Jeff Z Pan. Rethinking memory in ai: Taxonomy, operations, topics, and future directions. *arXiv preprint arXiv:2505.00675*, 2025.
- [65] Yaxiong Wu, Sheng Liang, Chen Zhang, Yichao Wang, Yongyue Zhang, Hufeng Guo, Ruiming Tang, and Yong Liu. From human memory to ai memory: A survey on memory mechanisms in the era of llms. *arXiv preprint arXiv:2504.15965*, 2025.
- [66] Yu Wang, Chi Han, Tongtong Wu, Xiaoxin He, Wangchunshu Zhou, Nafis Sadeq, Xiusi Chen, Zexue He, Wei Wang, Gholamreza Haffari, et al. Towards lifespan cognitive systems. *arXiv preprint arXiv:2409.13265*, 2024.
- [67] Tenghao Huang, Kinjal Basu, Ibrahim Abdelaziz, Pavan Kapanipathi, Jonathan May, and Muhan Chen. R2d2: Remembering, reflecting and dynamic decision making for web agents. *arXiv preprint arXiv:2501.12485*, 2025.
- [68] Rana Salama, Jason Cai, Michelle Yuan, Anna Currey, Monica Sunkara, Yi Zhang, and Yassine Benajiba. Meminsight: Autonomous memory augmentation for llm agents. *arXiv preprint arXiv:2503.21760*, 2025.
- [69] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *arXiv preprint arXiv:2005.11280*, 2020.
- [70] Larry R Squire, Lisa Genzel, John T Wixted, and Richard G Morris. Memory consolidation. *Cold Spring Harbor perspectives in biology*, 7(8):a021766, 2015.
- [71] Ye Ye. Task memory engine (tme): Enhancing state awareness for multi-step llm agent tasks. *arXiv preprint arXiv:2504.08525*, 2025.
- [72] Zhenhailong Wang, Haiyang Xu, Junyang Wang, Xi Zhang, Ming Yan, Ji Zhang, Fei Huang, and Heng Ji. Mobile-agent-e: Self-evolving mobile assistant for complex tasks. *arXiv preprint arXiv:2501.11733*, 2025.
- [73] Jiale Wei, Xiang Ying, Tao Gao, Fangyi Bao, Felix Tao, and Jingbo Shang. Ai-native memory 2.0: Second me. *arXiv preprint arXiv:2503.08102*, 2025.
- [74] Qitan Lv, Jie Wang, Hanzhu Chen, Bin Li, Yongdong Zhang, and Feng Wu. Coarse-to-fine highlighting: Reducing knowledge hallucination in large language models. *arXiv preprint arXiv:2410.15116*, 2024.
- [75] Kurt Shuster, Spencer Poff, Moya Chen, Douwe Kiela, and Jason Weston. Retrieval augmentation reduces hallucination in conversation. *arXiv preprint arXiv:2104.07567*, 2021.
- [76] Mengjia Niu, Hao Li, Jie Shi, Hamed Haddadi, and Fan Mo. Mitigating hallucinations in large language models via self-refinement-enhanced knowledge retrieval. *arXiv preprint arXiv:2405.06545*, 2024.
- [77] Chandana Sree Mala, Gizem Gezici, and Fosca Giannotti. Hybrid retrieval for hallucination mitigation in large language models: A comparative analysis. *arXiv preprint arXiv:2504.05324*, 2025.
- [78] Yao Xu, Shizhu He, Jiabei Chen, Zihao Wang, Yangqiu Song, Hanghang Tong, Guang Liu, Kang Liu, and Jun Zhao. Generate-on-graph: Treat llm as both agent and kg in incomplete knowledge graph question answering. *arXiv preprint arXiv:2404.14741*, 2024.

- [79] Yongrui Chen, Junhao He, Linbo Fu, Shenyu Zhang, Rihui Jin, Xinbang Dai, Jiaqi Li, Dehai Min, Nan Hu, Yuxin Zhang, et al. Pandora: A code-driven large language model agent for unified reasoning across diverse structured knowledge. *arXiv preprint arXiv:2504.12734*, 2025.
- [80] Minjae Seo, Wonwoo Choi, Myoungsung You, and Seungwon Shin. Autopatch: Multi-agent framework for patching real-world cve vulnerabilities. *arXiv preprint arXiv:2505.04195*, 2025.
- [81] Xiangru Tang, Tianyu Hu, Muyang Ye, Yanjun Shao, Xunjian Yin, Siru Ouyang, Wangchunshu Zhou, Pan Lu, Zhuosheng Zhang, Yilun Zhao, Arman Cohan, and Mark Gerstein. Chemagent: Self-updating library in large language models improves chemical reasoning. *arXiv preprint arXiv:2501.06590*, 2025.
- [82] Yusuf Roohani, Andrew Lee, Qian Huang, Jian Vora, Zachary Steinhart, Kexin Huang, Alexander Marson, Percy Liang, and Jure Leskovec. Predicting single-cell perturbation responses via graph neural networks with application to covid-19. *arXiv preprint arXiv:2405.17631*, 2024.
- [83] Kexin Chen, Junyou Li, Kunyi Wang, Yuyang Du, Jiahui Yu, Jiamin Lu, Lanqing Li, Jiezhong Qiu, Jianzhang Pan, Pheng Ann Heng, et al. Chemist-x: Large language model-empowered agent for reaction condition recommendation in chemical synthesis. *arXiv preprint arXiv:2311.10776*, 2023.
- [84] Jiayi Zhang, Jinyu Xiang, Zhaoyang Yu, Fengwei Teng, Xiong-Hui Chen, Jiaqi Chen, Mingchen Zhuge, Xin Cheng, Sirui Hong, Jinlin Wang, Bingnan Zheng, Bang Liu, Yuyu Luo, and Chenglin Wu. AFLOW: Automating agentic workflow generation. In *Proc. International Conference on Learning Representations (ICLR)*, 2025.
- [85] Longtao Zheng, Rundong Wang, Xinrun Wang, and Bo An. Synapse: Trajectory-as-exemplar prompting with memory for computer control. In *Proc. International Conference on Learning Representations (ICLR)*, 2024.
- [86] Jie Liu, Pan Zhou, Yingjun Du, Ah-Hwee Tan, Cees GM Snoek, Jan-Jakob Sonke, and Efstratios Gavves. Capo: Cooperative plan optimization for efficient embodied multi-agent cooperation. *arXiv preprint arXiv:2411.04679*, 2024.
- [87] Hanchao Liu, Rongjun Li, Weimin Xiong, Ziyu Zhou, and Wei Peng. Workteam: Constructing workflows from natural language with multi-agents. *arXiv preprint arXiv:2503.22473*, 2025.
- [88] Ruiwen Zhou, Yingxuan Yang, Muning Wen, Ying Wen, Wenhao Wang, Chunling Xi, Guoqiang Xu, Yong Yu, and Weinan Zhang. Trad: Enhancing llm agents with step-wise thought retrieval and aligned decision. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 3–13, 2024.
- [89] Mathis Pink, Qinyuan Wu, Vy Ai Vo, Javier Turek, Jianing Mu, Alexander Huth, and Mariya Toneva. Position: Episodic memory is the missing piece for long-term llm agents. *arXiv preprint arXiv:2502.06975*, 2025.
- [90] Minttu Alakuijala, Ya Gao, Georgy Ananov, Samuel Kaski, Pekka Marttinen, Alexander Ilin, and Harri Valpola. Memento no more: Coaching ai agents to master multiple tasks via hints internalization. *arXiv preprint arXiv:2502.01562*, 2025.

Limitations

Despite the promising empirical outcomes exhibited by AGENT KB, the proposed approach encounters intrinsic scalability constraints as AGENT KB undergoes expansion. The current retrieval mechanism, while demonstrating efficacy at the experimental scale, exhibits polynomial-time complexity scaling with respect to the number of stored experiential records. As the repository evolves from thousands to millions of entries spanning heterogeneous domains, maintaining sub-second retrieval latency becomes progressively intractable, thereby imposing potential limitations on real-time applications that demand immediate response capabilities.

The quality and reliability of automatically generated experiential knowledge represent a foundational limitation. Although extant validation mechanisms are capable of filtering overt failures, nuanced flaws within reasoning frameworks or domain-specific idiosyncrasies may still infiltrate the system undetected. More critically, the existing architecture remains devoid of a systematic experience deprecation mechanism or version control paradigm, such that obsolete or suboptimal strategies may persist indefinitely without periodic review. This structural deficiency not only impedes the system's adaptive evolutionary trajectory but also exacerbates the accumulation of suboptimal solutions in long-term operational scenarios.

Cross-domain knowledge transfer, while theoretically beneficial, demonstrates diminishing returns when applied to domains with minimal structural similarity. This phenomenon suggests fundamental boundaries to the universality of the proposed approach, necessitating careful consideration of domain-relatedness during knowledge base construction. Additionally, the reliance on pre-trained language models for experiential encoding and retrieval introduces an implicit bias toward tasks well-represented in the models' training corpora, potentially disadvantaging nascent or specialized domains.

Furthermore, the scope of domain knowledge embedded within the current system manifests significant incompleteness. Insufficient coverage of specialized domains engenders inherent limitations in addressing cross-disciplinary tasks, particularly in scenarios requiring deep integration of multi-field expertise. Such epistemic gaps further amplify the risk of reasoning errors in complex application environments, underscoring the urgency of establishing a dynamic knowledge updating mechanism and a comprehensive domain coverage framework.

Future Work

Advancing beyond retrieval-based knowledge reuse, we envision developing a causal reasoning framework that understands why certain strategies succeed in specific contexts. This framework would decompose experiences into causal chains, identifying prerequisite conditions, action-outcome relationships, and contextual dependencies. By modeling these causal structures explicitly, agents could synthesize novel solutions by recombining causal fragments rather than merely adapting complete experiences. Preliminary investigations suggest that causal decomposition could improve transfer effectiveness by 30-40% for cross-domain applications, particularly in scenarios requiring creative problem-solving rather than pattern matching.

The integration of continual learning mechanisms represents another crucial direction for AGENT KB’s evolution. Rather than treating AGENT KB as a static repository, we propose implementing experience refinement loops that automatically update strategies based on deployment outcomes. This would involve tracking the success rates of retrieved experiences in novel contexts, identifying systematic failure patterns, and synthesizing improved versions through automated experimentation. Such a system would require careful balance between exploration of new strategies and exploitation of proven approaches, potentially leveraging multi-armed bandit algorithms or evolutionary optimization techniques to guide the refinement process.

Theoretical foundations for cross-agent knowledge transfer remain underdeveloped, presenting opportunities for fundamental research. We plan to investigate formal frameworks for characterizing experience transferability, potentially drawing from domain adaptation theory and meta-learning. Understanding the geometric properties of experience embeddings and their relationship to task similarity could enable more principled retrieval mechanisms. Furthermore, developing provable guarantees for retrieval quality and transfer effectiveness would enhance AGENT KB’s applicability in high-stakes scenarios where performance bounds are critical.

Broad Impact

AGENT KB fundamentally transforms how AI systems accumulate and share knowledge, potentially accelerating the pace of AI development while reducing duplicated efforts across the research community. By enabling smaller organizations and individual researchers to leverage experiences accumulated by well-resourced institutions, our framework democratizes access to advanced problem-solving strategies. This democratization effect could be particularly transformative in developing countries and underfunded research areas, where limited computational resources currently constrain AI advancement. However, this concentration of knowledge also raises questions about intellectual property and competitive advantage, requiring careful consideration of contribution attribution and usage rights.

The transparency and interpretability afforded by AGENT KB's experience-based reasoning address growing concerns about AI accountability in critical applications. Unlike black-box neural systems, agents using AGENT KB can justify decisions by citing specific past experiences and the reasoning patterns derived from them. This traceability becomes invaluable in regulated industries such as healthcare and finance, where decision audit trails are legally mandated. Nevertheless, the system's reliance on historical experiences may inadvertently perpetuate past biases or outdated practices, notably if AGENT KB lacks diversity in contributors or problem domains.

The societal implications of widespread AGENT KB adoption extend beyond technical considerations. In educational settings, students can access expert problem-solving strategies that were previously available only through direct mentorship, potentially revolutionizing how complex skills are taught and learned. In professional contexts, AGENT KB could serve as an intelligent assistant that captures and propagates organizational knowledge, preventing expertise loss due to employee turnover. However, this same capability raises concerns about job displacement and the commoditization of expert knowledge. Ensuring that AGENT KB enhances rather than replaces human expertise requires thoughtful deployment strategies and ongoing dialogue between technologists, domain experts, and affected communities.

Acknowledgments

We would like to express our sincere gratitude to Xuanzhang Liu and Yigeng Jiang for their contributions and assistance during the revision process of this paper.

A Experience Representation and Storage

While the main paper focuses on our three core innovations (knowledge abstraction, dual-phase retrieval, and adaptive refinement), this appendix provides detailed technical descriptions of the mechanisms underlying experience abstraction, representation, and storage within AGENT KB.

A.1 Experience Abstraction

To ensure AGENT KB contains reusable and generalizable planning knowledge, we implement a structured abstraction pipeline that transforms raw execution logs from multiple agent frameworks into standardized experience patterns.

Raw Log Collection and Preprocessing. We begin by collecting execution logs from completed tasks across various domains and agent frameworks. These logs include:

- Input prompts and task descriptions
- Agent-generated reasoning steps
- Tool calls and their parameters
- Execution results and error messages

Each log is annotated with metadata including the agent framework used, domain type, and task difficulty level.

Before abstraction, logs undergo preprocessing to:

- Normalize tool names and parameter formats across frameworks
- Remove framework-specific implementation details
- Extract high-level reasoning traces and decision points

Manual Error Analysis and Pattern Identification. To identify common failure modes and improve generalization, human annotators manually inspect a subset of failed logs. They summarize recurring issues such as, *incorrect tool selection*, *misaligned reasoning chains*, or *missing preconditions or constraints*. These failures are abstracted into correction templates that serve as few-shot examples for the experience generation model.

Reasoning Template Design. The abstraction process relies on a set of reasoning templates , which define how experiences should be represented in AGENT KB. Each template includes slots for:

Agent KB data template
{ "question": "<question from various data source>", "agent_plan": "<Agent original plan>", "agent_experience": "<detailed agent experience>", } }

This structure ensures consistency across different frameworks and enables cross-task adaptation during retrieval.

Few-shot Prompting for Experience Generation. We construct task-specific prompts using manually designed reasoning templates. The complete set of prompts used in our experiments is detailed in Appendix O. Each prompt P comprises three core components:

$$P = \langle T, W, E_{\text{few-shot}} \rangle$$

where:

- T : the task description, which specifies input constraints and defines the objective
- W : the agent-generated workflow trace, a preprocessed and normalized representation of the solution path derived from raw execution logs

- $E_{\text{few-shot}}$: a set of k few-shot exemplars, each containing a previously abstracted experience entry encoded in a standardized format

Given this prompt, an LLM-based generator produces a structured experience entry E_{gen} :

$$E_{\text{gen}} = \text{LLM}(P) = \langle \pi, \gamma, S, \mathcal{C}, \mathcal{R} \rangle$$

Each generated experience follows the structure defined in Section 3.1. This formulation enables the systematic generation of high-quality, semantically consistent knowledge entries across diverse tasks and domains, facilitating scalable knowledge accumulation within AGENT KB.

A.2 Experience Representation

Each experience in AGENT KB is encoded as a structured tuple $E = \langle \pi, \gamma, S, \mathcal{C}, \mathcal{R} \rangle$, where:

- π represents the problem pattern, including task type, input structure, and constraints
- γ (Optional) denotes the goal or objective, including success criteria and expected outputs
- $S = \{s_1, s_2, \dots, s_n\}$ is a workflow capturing a sequence of reasoning and execution steps
- \mathcal{C} (Optional) captures contextual features including domain D and difficulty level δ
- \mathcal{R} a set of references to related experiences, representing relationships such as abstraction, composition, adaptation, and alternatives

This comprehensive representation enables AGENT KB to capture not only what worked but also contextual factors that influence success and alternative approaches that might be relevant in different scenarios.

Experience Representation and Organization Before detailing the retrieval process, we define how experiences are represented within AGENT KB. Each experience E is encoded with multi-faceted embeddings:

$$f(E) = \{f^\pi, f^\gamma, f^S, f^C\},$$

where f^π represents the problem pattern embedding, f^γ the goal embedding, f^S the solution steps embedding, and f^C the context embedding (optional).

Experiences are organized in a hierarchical knowledge graph $\mathcal{KB} = (V, \mathcal{E})$, where nodes V represent individual experiences and edges \mathcal{E} encode semantic and structural relationships between them (e.g., abstraction, composition, alternatives). While the knowledge graph structure is currently used only for experience organization and visualization, it lays the foundation for future extensions that may incorporate relational reasoning or graph-based retrieval mechanisms into the experience reuse pipeline.

Experiences are organized in a hierarchical knowledge graph. While this graph primarily supports navigation and visualization, it also provides a foundation for future extensions that may incorporate relational reasoning into retrieval.

Student Agent: Query-based Workflow Retrieval When a query Q (e.g., a GAIA benchmark problem) is received, the student agent initiates the first retrieval phase. The student agent first *reasons* about how to approach the problem, identifying key requirements and potential solution strategies. Then, it performs *retrieval* from AGENT KB to find relevant experiences that might guide its planning process. Given the current agent state \hat{S} with problem $\hat{\pi}$ and goal $\hat{\gamma}$, the student retrieves relevant experiences through:

$$\text{Retrieve}(\hat{S}, \hat{\pi}, \hat{\gamma}, k) = \arg \top_k \left(\text{sim}(E_i, \hat{S}) \cdot \text{Relevance}(E_i, \hat{\pi}, \hat{\gamma}) \cdot \text{Success}(E_i) \right)$$

This query-based retrieval process operates through a sophisticated multi-stage approach that balances broad similarity matching with precise state alignment.

First, we perform coarse retrieval based on problem-goal similarity, identifying experiences where:

$$\mathcal{E}_{\text{coarse}} = \{E_i \mid \text{sim}_{\text{cos}}(r_i, \hat{r}) > \theta_{\text{coarse}}\}$$

where $r_i = f(\pi_i) + f(\gamma_i)$ and $\hat{r} = f(\hat{\pi}) + f(\hat{\gamma})$. This is complemented by fine-grained retrieval that uses the current agent state to find experiences with matching execution steps, where:

$$\mathcal{E}_{\text{fine}} = \arg \top_k(S_i^{\text{fine}})$$

with:

$$S_i^{\text{fine}} = \sum_{j=1}^m \max_{\ell=1, \dots, L_i} \text{sim}_{\text{cos}}(s_{i,\ell}, \hat{s}_j)$$

Finally, these retrieval strategies are combined through an adaptive mechanism:

$$S_i(t) = \lambda(t) \cdot S_i^{\text{coarse}} + (1 - \lambda(t)) \cdot S_i^{\text{fine}}$$

where $S_i^{\text{coarse}} = \text{sim}_{\text{cos}}(r_i, \hat{r})$ and $\lambda(t) \in [0, 1]$ is a time-dependent weighting function that balances coarse and fine-grained retrieval based on the current stage of problem solving.

This creates a context-sensitive retrieval approach that evolves throughout the problem-solving process, with final selection given by:

$$\mathcal{R}(t) = \arg \top_k S_i(t)$$

The retrieved experiences contain successful workflows from similar historical tasks, including critical elements such as complete planning structures (step sequences), appropriate tool selection for each step, and general reasoning patterns relevant to the query type. The student agent's primary focus at this stage is ensuring the overall workflow structure is appropriate for the task.

The student agent then adapts these experiences to the current context, applying operations such as parameter substitution, step expansion/contraction, and domain translation:

$$E_{\text{adapted}} = \text{Adapt}(E_{\text{retrieved}}, \hat{S}, \hat{\pi}, \hat{\gamma})$$

These adapted experiences are synthesized to generate an initial execution plan:

$$\text{Plan}_{\text{initial}} = \text{Integrate}(\text{Plan}_{\text{empty}}, \{E_{\text{adapted}}\})$$

This plan includes a sequence of reasoning steps $S = \{s_1, s_2, \dots, s_n\}$, each with specified tools and execution parameters. The student agent executes this plan, generating execution logs L that capture both successes and failures during the process.

Teacher Agent: Log-based Reasoning and Refinement After the initial execution, the student agent forwards both the query Q and execution logs L to the teacher agent. Unlike the student agent, which focuses on planning the overall workflow, the teacher agent performs critical reasoning functions on the execution itself. The teacher agent analyzes the logs through three main processes:

First, it performs error analysis to identify problematic steps:

$$\mathcal{E}(L) = \{(s_i, \text{error}_i, \text{cause}_i) \mid s_i \in S, \mathbb{H}(s_i) = \text{True}\},$$

where \mathcal{E} denotes the error analysis function, extracting error steps with causes. Next, it summarizes the execution log to extract key patterns:

$$\mathcal{L}(L) = \mathbb{S}(\{s_1, s_2, \dots, s_n\}),$$

where \mathcal{L} represents log summarization, abstracting execution patterns. Finally, it evaluates the overall performance by comparing actual outcomes with expected results:

$$\mathcal{P}(L, Q) = \mathbb{E}(\mathcal{O}(L), \mathcal{O}^*(Q)),$$

where \mathcal{P} performs performance evaluation, comparing actual $\mathcal{O}(L)$ and expected $\mathcal{O}^*(Q)$ outcomes. Based on this comprehensive analysis, the teacher agent identifies problematic steps that require refinement:

$$\mathcal{S}_p = \mathcal{I}(\mathcal{E}(L), \mathcal{P}(L, Q)),$$

where \mathcal{S}_p is the set of steps requiring refinement, derived from error analysis and performance evaluation via \mathcal{I} (issue identification function).

For each problematic step s_i , the teacher agent performs a targeted secondary retrieval from AGENT KB, focusing on fine-grained matching of similar step-level experiences:

$$E_{\text{refinement}} = \arg \top_m \left(\sum_{s_i \in \mathcal{S}_p} \max_l \mathcal{C}(s_i, s_l^{(j)}) \cdot \text{Precision}(E_j) \right)$$

Here, \mathcal{C} denotes cosine similarity function; $\text{Precision}(E_j)$ is the precision metric for experience E_j .

Unlike the first retrieval phase, which focused on high-level workflow structure, this log-based refinement retrieval targets specific execution details that affect precision and correctness. The teacher agent identifies granular aspects such as precise parameter configurations (e.g., maintaining three decimal places in calculations), error handling strategies for specific failure modes, tool usage refinements and constraints, and step-specific reasoning patterns that improve accuracy. These fine-grained execution details are critical for successfully completing tasks that require not just the right approach but also precise implementation.

The teacher agent then adapts these refinement experiences: $E_{\text{refined}} = \mathcal{A}(E_{\text{refinement}}, L, Q)$, where \mathcal{A} is the experience adaptation function, integrating logs and queries into retrieved experiences. Based on the refined experiences, the agent generates a set of specific refinement hints through reasoning:

$$\mathcal{H} = \mathcal{G}(E_{\text{refined}}, \mathcal{S}_p, L, Q),$$

where \mathcal{G} generates actionable hints by reasoning over adapted experiences and problematic steps \mathcal{S}_p .

Benefits of the Dual-Phase Approach This two-phase approach significantly enhances performance by addressing both structural correctness and execution precision. The Query-based Retrieval ensures the overall workflow structure is appropriate for the task (correct sequence of steps and tool selection), while the Log-based Refinement focuses on execution details that impact success (precise calculations, error handling, parameter tuning).

Through this teacher-student collaboration, AGENT KB enables progressive refinement that mimics human expert-apprentice learning relationships. Both agents employ the **Reason-Retrieve-Refine** pipeline, but with different focuses: the student agent reasons about the problem structure and overall solution approach, while the teacher agent reasons about the execution quality and potential improvements. The teacher agent effectively transfers knowledge from past experiences to guide the student agent toward successful task completion, with each phase targeting a different aspect of performance improvement.

A.3 Storage and Indexing

Experiences are organized in a hierarchical knowledge graph $\mathcal{KB} = (V, \mathcal{E})$ where vertices V represent individual experiences and edges \mathcal{E} denote meaningful relationships between them. These relationships include:

- **Abstraction:** connecting concrete experiences to their abstracted versions
- **Composition:** linking sub-workflows to composite workflows
- **Adaptation:** connecting experiences that have been successfully adapted across domains
- **Alternative:** connecting different approaches to solving similar problems

This graph structure facilitates efficient navigation across related experiences, enabling both breadth-first exploration of alternatives and depth-first exploration of hierarchical solution approaches.

To enable efficient retrieval over this structured repository, we employ a multi-indexing strategy. Specifically, two primary indices form the basis of the retrieval mechanism:

- **Semantic index:** Encodes the semantic meaning of problems and goals to enable intent-driven retrieval, identifying experiences addressing conceptually similar tasks.
- **Structural index:** Captures workflow structure patterns to support retrieval based on similarities in process organization or control flow.

Together, these indices underpin a dual-phase retrieval approach that efficiently identifies relevant experiences at both the workflow and component level, avoiding exhaustive traversal of the entire knowledge graph.

B Detailed Retrieval Mechanisms

B.1 Agent Retrieval Process

Both the STUDENT AGENT and TEACHER AGENT employ a dual-phase retrieval mechanism over AGENT KB to identify relevant past experiences. The retrieval process is algorithmically identical for both agents, with the distinction lying in the target of retrieval: the Student Agent retrieves workflows that match its current task, while the Teacher Agent retrieves guidance strategies or interventions aligned with the student’s needs.

B.1.1 Workflow Retrieval Function

The primary retrieval function for both agents is defined as:

$$\mathcal{E}_{\text{retrieved}} = \arg \text{top}_k(S_{\text{workflow}}(E_i, T))$$

where T denotes the retrieval target — for the Student Agent, this corresponds to the query or reasoning trace of the current task; for the Teacher Agent, it reflects the combination of execution log conclusions and the initial problem-solving plan.

The similarity score S_{workflow} is computed as a weighted sum of cosine similarities across multiple components (policy, context, etc.), using their respective embedding representations.

B.1.2 Agent-Specific Retrieval Targets and Sample Matching

The AGENT KB stores structured experiences, each encoded as a tuple:

$$E = (\pi, \gamma, S, \mathcal{C}, \mathcal{R})$$

Each agent retrieves experiences by matching its retrieval target against different components of the stored experience:

- **Student Agent:** Retrieves experiences by matching the current task’s query or reasoning trace against the π (and optionally γ) of stored experiences. This corresponds to retrieving semantically similar tasks based on problem patterns and goals.
- **Teacher Agent:** Retrieves experiences by matching observed execution logs and student behavior against the S (workflow) component of stored experiences. This allows it to find pedagogically relevant strategies that share similar reasoning paths or solution structures.

In both cases, the whole experience E from the matched sample is returned for further adaptation or analysis. This shared retrieval mechanism enables both agents to benefit from a unified knowledge representation while serving distinct functional roles.

C Adaptation Functions for Cross-Domain Transfer

This appendix provides additional technical details on the adaptation functions used for cross-domain knowledge transfer in AGENT KB.

C.1 Workflow-Level Adaptation

The workflow-level adaptation function transforms retrieved workflows to match the target domain's characteristics:

$$P_{\text{adapted}} = \mathcal{F}_{\text{adapt}}(P_{\text{retrieved}}, C_{\text{source}}, C_{\text{target}})$$

where $\mathcal{F}_{\text{adapt}}$ denotes the contextual adaptation operator that applies a series of domain-aware transformations. This function performs the following key operations:

- **Entity mapping:** replaces domain-specific entities with their target domain equivalents
- **Tool substitution:** replaces tools used in the source domain with equivalent tools in the target domain
- **Step reordering:** adjusts the sequence of steps to account for different dependencies in the target domain
- **Constraint handling:** modifies steps to accommodate different constraints in the target domain

The mapping function takes domain representations as input and outputs transformation parameters:

$$\theta_{\text{transform}} = f_{\text{map}}(C_{\text{source}}, C_{\text{target}})$$

These parameters are then used to guide the specific transformations applied to each component of the workflow.

C.2 LLM-Based Reasoning and Experience Refinement

To enhance cross-domain adaptation practically and efficiently, we propose an approach that combines large language model (LLM)-based reasoning with iterative experience refinement. Instead of relying on predefined transformation rules or dedicated neural modules to learn domain mappings, we prompt a pre-trained LLM with structured representations of the source and target domains, C_{source} and C_{target} , enabling it to infer and apply appropriate domain-specific transformations directly to the workflow.

The LLM infers transformation parameters from domain representations:

$$\theta_{\text{transform}} = \text{LLM}(C_{\text{source}}, C_{\text{target}})$$

These parameters guide the adaptation of each workflow component:

$$P_{\text{adapted}} = \text{Transform}(P_{\text{retrieved}}, \theta_{\text{transform}})$$

To refine adaptation over time, successful transformations are stored in a feedback pool and reused as exemplars in few-shot prompts:

$$(C_{\text{source}}, C_{\text{target}}, P_{\text{adapted}}) \in E_{\text{feedback}} \quad \text{if adaptation succeeds}$$

This iterative refinement enables the system to accumulate cross-domain adaptation experience, thereby improving future transfers through more informed LLM prompting.

D Experimental Details

D.1 Experimental Cost

All services used in this work rely on third-party API calls to OpenAI’s language models (GPT-4.1, Claude-3-7-sonnet, o1, etc.). The total cost of execution is primarily determined by the number of tokens processed during both prompt input and model output generation. Specifically, we report the token cost associated with different modules of our AGENT KB (Knowledge Base) system, as well as the per-agent token consumption during task execution.

As summarized in Table 6, the token cost of GPT-4.1 varies significantly depending on the complexity of the agent and its interaction with AGENT KB. For instance, the Base Model requires a relatively high number of reasoning steps (up to 12), resulting in a higher cumulative token count across multiple interactions. In contrast, the Student Agent and Teacher Agent, while still utilizing LLM-based inference, operate in a more passive or structured manner, resulting in fewer dynamic interactions and correspondingly lower token usage. The Database Generation module incurs a one-time cost during initialization, where large volumes of domain-specific knowledge are encoded into structured prompts for retrieval-augmented generation.

Given that OpenAI pricing is typically calculated based on both input and output tokens, the total cost of our experiments remains moderate due to suitable prompt engineering and step-limited execution strategies.

Table 6: Analysis of computational costs on the GAIA benchmark for AGENT KB. All costs, excluding database generation, correspond to a single evaluation on the GAIA validation set (165 tasks).

Type	Module	Prompt Tokens	Completion Tokens	Cost	Max Steps
Base Model	Action	~34M	~7M	~\$84.32	12
Database Generation		~5M	~750K	~\$10.88	-
Log summary	AGENT KB	~1M	~10K	~\$1.41	-
Student agent		~35K	~15K	~\$0.13	-
Teacher agent		~45K	~15K	~\$0.14	-

Token prices: \$1.36/M prompt token, \$5.44/M completion token.

As shown in Table 7, the computational costs of SWE-bench evaluation under the AGENT KB framework vary based on the source and structure of the hint material. Reasoning modules using RepoClassBench incur higher token costs due to deeper reasoning chains and longer hint contexts. In contrast, lightweight configurations, such as Top-n SWE-Gym with shorter hints and fewer reasoning steps, significantly reduce the per-item cost. By tailoring the prompt size and controlling the number of refinement steps, we maintain a low average cost (under \$0.008 per instance), ensuring the framework is scalable for large-scale software engineering benchmarks.

Table 7: Analysis of computational costs on the SWE-bench benchmark for AGENT KB modules. All costs correspond to per-item inference using GPT-4.1

Hint Source	Module	Prompt Tokens	Completion Tokens	Cost (/item)	Hint Length (tokens/item)	Max Steps
RepoClassBench	Reasoning	~6.5K	~850	~\$0.007805	~90	100
RepoClassBench	Refine	~4.2K	~450	~\$0.0028	~130	100
Top-n SWE-Gym	Retrieval+Refine	~2.8K	~300	~\$0.001875	~60	100
Top-n RepoClassBench	Retrieval+Refine	~3.1K	~350	~\$0.002125	~70	100

Token prices: \$1.36/M prompt tokens, \$5.44/M completion tokens.

D.2 Ablation Study on Knowledge Base Size

To evaluate how knowledge base size impacts retrieval and task-solving performance, we conduct ablation experiments using two KB configurations: *KB-100* (100 randomly sampled entries) and *KB-500* (500 randomly sampled entries).

Table 8: Performance of AGENT KB with varying knowledge base sizes on GAIA and SWE-bench.

Size	GAIA			SWE-bench	
	Average	Level 1	Level 2	Level 3	Resolved
Baseline	61.21	79.25	58.14	34.62	28.33
<i>KB-100</i>	56.97	73.58	53.49	34.62	24.67
<i>KB-500</i>	58.79	77.36	54.65	34.62	25.33

From Table 8, the results show a consistent trend: as AGENT KB size increases, overall performance improves. The full knowledge base achieves the highest scores across both benchmarks, indicating that access to more comprehensive information enhances the agent’s ability to retrieve relevant knowledge and solve complex tasks.

On GAIA, reducing the KB size leads to a gradual decline in average performance, from 61.21% to 58.79% with *KB-500* and further to 56.97% with *KB-100*. Notably, performance drops most significantly at Level 1 and Level 2 tasks, suggesting that larger KBs better support basic factual retrieval and moderate reasoning. Level 3 performance remains constant across configurations, implying that very challenging tasks may not benefit substantially from increased KB size due to other limiting factors such as model capacity or task complexity.

Similarly, on SWE-bench, the full KB achieves the best result (28.33% resolved), followed by *KB-500* and *KB-100*. This indicates that even for code-centric problem-solving, access to a broader knowledge base contributes positively to performance.

In summary, increasing knowledge base size generally enhances retrieval and task-solving capabilities, especially for less complex tasks. However, the diminishing returns observed suggest that there may be a point of saturation where additional knowledge yields marginal gains..

D.3 Ablation Details of Reason-Retrieve-Refine Modules

To evaluate the effectiveness of each component in our AGENT KB framework, we conduct a series of ablation studies. Our system consists of two agents: the Student Agent and the Teacher Agent, with distinct roles across two reasoning stages.

- **Student Agent** is responsible for the initial stage, which begins with **Reason** (to summarize key features from the input), followed by **Retrieve** (to find relevant prior experiences), and concludes with **Refine** (to improve the suggestions based on retrieved information).
- **Teacher Agent** operates in the second stage, where it begins with **Reason** (to analyze the logs and identify key errors), followed by **Retrieve** (to gather relevant experience), and concludes with **Refine** (to improve or correct the suggestions based on the retrieved information).

The experimental setup involves systematically removing or disabling specific modules or agents to assess their contributions. The results are summarized in Table 3, with the following definitions:

- w/o Student Agent: The first-stage steps are removed.
- w/o Teacher Agent: The second-stage steps are removed.
- w/o **Reason** Module : In both stages, no reasoning is performed; only retrieval based on raw data is conducted.
- w/o **Retrieve** Module : Both stages omit the retrieval process entirely. Agents rely solely on prompt-based instructions to generate responses, without consulting prior experiences.
- w/o **Refine** Module : No refinement is performed of both stags; only the retrieved content is used as knowledge.
- w/ Raw Workflow : The full retrieve pipeline is used, but without any explicit modular control—i.e., the model follows a standard prompting strategy throughout, lacking structured guidance through the **Reason** and **Refine** phases.

These ablation experiments provide insight into how each module contributes to overall performance, particularly in terms of accuracy, robustness, and coherence in complex reasoning tasks.

D.4 GAIA Details

Evaluated on the validation set of GAIA across three difficulty levels:

- **Level 1 (53 tasks):** Basic tasks requiring simple reasoning or straightforward retrieval.
- **Level 2 (86 tasks):** Intermediate complexity with multi-step reasoning or tool usage.
- **Level 3 (26 tasks):** Advanced tasks demanding sophisticated reasoning and domain knowledge.

Performance is measured using an unweighted average over all 165 tasks.

Two metrics are used:

- Pass@1: Evaluates the correctness of the first generated solution.
- Pass@3: Evaluates whether any of the three independently generated solutions is correct.

Method Configurations:

- `+AGENT KB / +AGENT KB ✓`: Evaluated using Pass@1, representing the model's initial attempt or after one round of feedback.
- `+AGENT KB ✓✓`: Uses Pass@3 to align with standard practices and improve comparability with existing methods.

It is essential to note that all experimental outcomes presented in this study are exclusively derived from evaluations conducted on the GAIA validation set rather than the test set. This methodological choice stems from two principal motivations. The first rationale pertains to comparative consistency - most established frameworks and baseline methodologies documented in prior research predominantly utilize validation set metrics, thereby enabling equitable performance comparisons across different systems. Second, our analytical framework, particularly concerning the implementation of `+AGENT KB ✓` and `+AGENT KB ✓✓` architectures, necessitates feedback mechanisms contingent upon the proper answer. This operational requirement mandates persistent access to verified ground-truth annotations for both incremental assessment and procedural optimization during inference phases. Since test set annotations remain conventionally undisclosed, we systematically employ the validation set as the foundation for all empirical analyses reported herein.

D.5 SWE-bench Details.

Performance is measured using an unweighted average over all 300 tasks.

Two metrics are used:

- Pass@1: Evaluates the correctness of the first generated solution.
- Pass@3: Evaluates whether any of the three independently generated solutions is correct.

Model Configurations:

- `+AGENT KB / +AGENT KB ✓`: Evaluated using Pass@1, representing the model's initial attempt or after one round of feedback.
- `+AGENT KB ✓✓`: Uses Pass@3 to align with standard practices and improve comparability with existing methods.

E Additional Details of Methodology

E.1 Experience Integration and Conflict Resolution

The teacher agent returns these refinement hints to the student agent, which must integrate them with the initial plan. This integration process requires resolving potential conflicts: $\text{Plan}_{\text{refined}} = \text{Integrate}(\text{Plan}_{\text{initial}}, \{\text{Hints}\})$, with conflict resolution following:

$$\text{Conflict}(p_1, p_2) = \begin{cases} \text{Merge}(p_1, p_2) & \text{if } \text{Compatible}(p_1, p_2) > \theta_c \\ \text{Select}(p_1, p_2) & \text{otherwise} \end{cases}.$$

The student agent then executes this refined plan, typically achieving superior performance compared to the initial execution.

E.2 Knowledge Evolution

AGENT KB continuously evolves through collaborative experience refinement: $E_{\text{refined}} = \text{Refine}(E, \mathcal{U})$, where \mathcal{U} is the usage history containing information about when and how the experience has been used. Similar experiences from different agents are merged:

$$E_{\text{merged}} = \text{Merge}(E_i, E_j) = \langle \pi_{ij}, \gamma_{ij}, S_{ij}, C_{ij}, \mu_{ij}, \mathcal{F}_{ij}, \mathcal{R}_{ij} \rangle,$$

while outdated or low-value experiences are pruned:

$$\text{Prune}(\mathcal{KB}) = \{E \in \mathcal{KB} | \text{Utility}(E, t_{\text{current}}) > \theta_p\},$$

with utility decaying over time unless reinforced:

$$\text{Utility}(E, t) = Q(E) \cdot e^{-\lambda(t-t_{\text{recent}})} + \sum_{i=1}^n \text{UsageImpact}(E, t_i),$$

The complete **Reason-Retrieve-Refine** pipeline operates within both the student and teacher agents, though with different objectives and contexts:

$$\text{RRR}(\hat{S}, \hat{\pi}, \hat{\gamma}) = \text{Refine}(\text{Retrieve}(\hat{S}, \hat{\pi}, \hat{\gamma}), \hat{S}),$$

and AGENT KB evolves according to:

$$\mathcal{KB}_{t+1} = \text{Update}(\mathcal{KB}_t, \{\text{Reason}(W_i)\}_{i=1}^{N_W}, \{\text{Feedback}(E_j)\}_{j=1}^{N_E}).$$

The framework-agnostic design enables different agents to both contribute to and benefit from the shared knowledge base, creating a virtuous cycle of collective intelligence improvement that enhances the performance of multi-agent systems over time.

E.3 Prospective Framework for Adaptive Experience Quality Refinement

The current methodological framework does not implement dynamic quality updates for experiential knowledge based on operational outcomes. To address this limitation, we outline a strategic research direction involving the development of an adaptive experience valuation mechanism that iteratively adjusts epistemic weights based on task-solving efficacy. Specifically, upon completion of each operational iteration, the quality parameter Q associated with deployed experience E would undergo recalibration through the following update rule:

$$Q_{\text{new}}(E) = (1 - \alpha) \cdot Q_{\text{old}}(E) + \alpha \cdot \mathcal{F}(E, \hat{S})$$

where $\alpha \in [0, 1]$ represents a tunable temporal discount factor modulating update intensity. At the same time, $\mathcal{F}(E, \hat{S})$ quantifies context-specific utility through a multidimensional assessment of solution \hat{S} 's task completion success. This meta-learning paradigm would enable progressive optimization of experiential hierarchies through reward attribution, enhancing both retrieval relevance through quality-aware prioritization and operational robustness in subsequent task executions. The proposed mechanism establishes formal foundations for autonomous competence refinement in evolving problem spaces.

F Retrieval Details

F.1 Retrieval Architecture

AGENT KB employs a two-stage retrieval framework designed to progressively refine the selection of relevant past experiences for effective task planning and execution:

Summary-based Retrieval. The second retrieval phase conducts a fine-grained analysis of execution logs (e.g., `intermediate_steps`) associated with the retrieved experiences. Specifically, we summarize both the overall plan structure and individual reasoning or action steps from these logs. These summaries are then used to perform a more detailed retrieval, aligning the current task state with specific subroutines or decision points from past executions. This step facilitates the identification of practical low-level actions or reasoning patterns that are contextually aligned with the current execution trajectory.

Criticism-Based Retrieval. The system actively searches for past experiences based on shared error patterns rather than task goals or outcomes. This stage focuses on identifying historical execution logs that contain similar types of mistakes—such as flawed reasoning steps, incorrect actions, or strategic misjudgments—as the current task. By encoding and matching these failure modes semantically, the retrieval process surfaces relevant cases where similar problems arose, allowing the planner to learn from prior failures and avoid repeating them. This error-driven approach enables a more proactive and reflective planning process grounded in lessons from past critiques.

F.2 Retrieval Types.

To ensure robust and contextually relevant experience retrieval, we incorporate multiple retrieval mechanisms that operate at different levels of abstraction. Within this framework, we utilize three primary types of retrieval: Text similarity retrieval, semantic retrieval, and hybrid retrieval, each offering distinct advantages in capturing relevance between the current task and historical experiences.

Text similarity retrieval. Text similarity retrieval is based on surface-level term matching and relies on traditional information retrieval techniques such as TF-IDF (Term Frequency-Inverse Document Frequency). This method quantifies the importance of terms within a document relative to a corpus, representing textual content as sparse, high-dimensional vectors. It excels at identifying documents that share significant keyword overlap with the query, making it particularly effective when vocabulary alignment is strong.

Semantic Retrieval. Semantic retrieval goes beyond keyword matching by encoding text into dense vector representations that capture meaning and contextual relationships. In our implementation, we use the `sentence-transformers/all-MiniLM-L6-v2` model, a lightweight yet powerful transformer-based encoder that maps sentences and paragraphs into a continuous vector space. This allows for the computation of cosine similarity between embeddings, enabling the system to retrieve experiences that are semantically related—even if they do not share exact text similarity overlap.

Hybrid Retrieval. To combine the strengths of both text similarity and semantic approaches, we also implement hybrid retrieval, which fuses results from both retrieval methods using a weighted ranking strategy. For a retrieved experience $e_i \in \mathcal{E}$, the final relevance score is computed as a linear combination of its text-based similarity score σ_i^{text} and semantic similarity score σ_i^{sem} :

$$\sigma_i^{\text{hyb}} = \alpha \cdot \sigma_i^{\text{text}} + (1 - \alpha) \cdot \sigma_i^{\text{sem}},$$

where $\alpha \in [0, 1]$ is a tunable parameter (default: $\alpha = 0.5$) that balances the influence of each retrieval modality. Hybrid retrieval offers a balanced trade-off between precision and generalization, mitigating the limitations of individual methods. It ensures that the retrieval mechanism remains robust to both syntactic variation and conceptual drift while maintaining interpretability and performance.

G Additional Experiment

G.1 Additional Evaluations

This section provides comprehensive results for the experiments conducted in the main text. We present detailed performance metrics across different models and retrieval strategies on the GAIA and SWE-bench, as well as ablation studies to analyze the effectiveness of our proposed components.

Table 9 presents the detailed performance of various large language models, including GPT-4o, GPT-4.1, o3-mini, Claude-3.7, Qwen-3 32B, and DeepSeek-R1, under different experimental settings. The evaluation includes baseline performance and improvements achieved by incorporating the +AGENT KB, +AGENT KB ✓, and +AGENT KB ✓♡ methods. Performance is measured using average accuracy and per-level accuracy on GAIA validation set, along with SWE-bench resolved scores. The final row (“Gap”) indicates the improvement from the baseline to the best-performing method for each model. Notably, all models show significant gains when using the enhanced reasoning and retrieval capabilities introduced by our framework.

Table 9: Detailed results of various base models on GAIA.

Model	Method	Average	GAIA Level 1	GAIA Level 2	GAIA Level 3	SWE-bench Resolved
GPT-4o	Baseline	45.06	62.26	45.35	15.38	16.33
	+AGENT KB	46.67	66.04	44.19	15.38	20.33
	+AGENT KB ✓	55.15	71.70	48.84	42.31	29.33
	+AGENT KB ✓♡	58.79	77.36	52.33	42.31	31.33
	Gap	$\Delta 13.73$	$\Delta 15.10$	$\Delta 6.98$	$\Delta 26.93$	$\Delta 15.00$
GPT-4.1	Baseline	55.15	67.92	53.49	34.62	24.33
	+AGENT KB	61.21	79.25	58.14	34.62	28.33
	+AGENT KB ✓	67.27	83.02	67.44	34.62	37.33
	+AGENT KB ✓♡	73.94	84.91	73.26	53.85	38.00
	Gap	$\Delta 18.79$	$\Delta 16.99$	$\Delta 19.77$	$\Delta 19.23$	$\Delta 13.67$
o3-mini	Baseline	32.12	47.17	26.74	19.23	23.00
	+AGENT KB	29.09	39.62	25.58	19.23	31.67
	+AGENT KB ✓	33.33	45.28	30.23	19.23	35.33
	+AGENT KB ✓♡	40.60	52.83	38.37	23.08	37.00
	Gap	$\Delta 8.48$	$\Delta 5.66$	$\Delta 11.63$	$\Delta 3.85$	$\Delta 14.00$
Claude-3.7	Baseline	58.79	64.15	61.63	38.46	30.00
	+AGENT KB	65.45	75.47	66.28	38.46	46.67
	+AGENT KB ✓	69.70	79.25	69.77	50.00	49.67
	+AGENT KB ✓♡	75.15	84.91	74.42	57.69	51.00
	Gap	$\Delta 16.36$	$\Delta 20.76$	$\Delta 12.79$	$\Delta 19.23$	$\Delta 9.67$
Qwen-3 32B	Baseline	35.76	47.17	38.37	3.85	18.33
	+AGENT KB	41.82	64.15	33.72	23.08	20.67
	+AGENT KB ✓	46.67	71.70	37.21	26.92	28.67
	+AGENT KB ✓♡	49.70	75.47	40.70	26.92	30.33
	Gap	$\Delta 13.94$	$\Delta 38.30$	$\Delta 2.33$	$\Delta 23.07$	$\Delta 12.00$
DeepSeek-R1	Baseline	49.70	62.26	50.00	23.08	24.33
	+AGENT KB	50.91	69.81	50.00	15.38	26.67
	+AGENT KB ✓	58.18	73.58	56.98	30.77	31.00
	+AGENT KB ✓♡	63.64	79.25	61.63	38.46	32.67
	Gap	$\Delta 13.94$	$\Delta 16.99$	$\Delta 11.63$	$\Delta 15.38$	$\Delta 8.34$

G.2 Retrieval Analysis

Table 10 compares summary-based and criticism-based retrieval methods across text similarity, semantic similarity, and hybrid strategies on GAIA and SWE-bench. Three key patterns emerge: (1) Hybrid retrieval achieves peak performance for summary-based methods (67.27 average on GAIA), while criticism-based methods perform best with text similarity (66.06 average). (2) Task complexity inversely correlates with performance across all methods, with Level 3 GAIA scores declining to 34.62-38.46% versus 73.58-83.02% for Level 1. (3) SWE-bench results show narrower margins between methods (4% resolved scores), suggesting benchmark-specific sensitivity to retrieval approaches.

The ablation study in Table 11 reveals three parameterization insights: (1) Optimal top-k values differ by method - text similarity peaks at k=3 (64.24 GAIA average), semantic similarity at k=5 (62.42), and hybrid search at k=3 (67.27). (2) Level 3 performance shows counterintuitive trends, with text similarity declining 7.7% from k=1 to k=5 while hybrid search improves 11.5%. (3) Parameter sensitivity varies substantially, with hybrid retrieval showing minimal $k = 1$ to $k = 5$ variance versus text similarity's 3.4% drop.

Cross-analysis identifies two critical interactions: (1) Summary-based hybrid retrieval with $k=3$ configuration achieves maximum GAIA performance (83.02% Level 1, 67.44% Level 2). (2) Criticism-based text similarity with $k=1$ yields best Level 3 results (38.46%), outperforming all hybrid configurations. These findings demonstrate that optimal retrieval configurations depend on both content type (summary vs. criticism) and task complexity, necessitating adaptive strategy selection rather than universal solutions.

Table 10: Retrieval results by different retrieval types on GAIA and SWE-bench.

Retrieval	Type	GAIA			SWE-bench Resolved
		Average	Level 1	Level 2	
Summary-based	Text Similarity	64.24	77.36	65.11	34.62
	Semantic similarity	58.79	69.81	59.30	34.62
	Hybrid search	67.27	83.02	67.44	34.62
Criticism-based	Text similarity	66.06	77.36	67.44	38.46
	Semantic similarity	62.42	73.58	63.95	34.62
	Hybrid search	63.03	77.36	62.79	34.62

Table 11: Retrieval performance across different top- k on GAIA and SWE-bench.

Retrieval Type	Top- k	GAIA			SWE-bench Resolved
		Average	Level 1	Level 2	
Text sim.	$k = 1$	63.03	75.47	62.79	38.46
	$k = 3$	64.24	77.36	65.11	34.62
	$k = 5$	62.42	77.36	62.79	30.77
Semantic sim.	$k = 1$	60.00	73.58	58.13	38.46
	$k = 3$	58.79	69.81	59.30	34.62
	$k = 5$	62.42	75.47	61.63	38.46
Hybrid.	$k = 1$	63.64	79.25	62.79	34.62
	$k = 3$	67.27	83.02	67.44	34.62
	$k = 5$	66.67	81.13	66.28	38.46

G.3 Knowledge Source Comparison

We also investigate the impact of different knowledge sources on AGENT KB performance. Table 12 compares performance using knowledge derived from different sources: Hand (manually crafted knowledge entries created by domain experts) and Generate (automatically generated knowledge entries derived from agent interactions). Additionally, we compare our method against SOTA (state-of-the-art results achieved by current closed-source agent frameworks on GAIA) and Open Source

(state-of-the-art results achieved by current open-source agent frameworks on GAIA). Interestingly,

Table 12: Performance comparison across different experience types on GAIA and SWE-bench.

Experience type	Average	GAIA			SWE-bench
		Level 1	Level 2	Level 3	Resolved
Hand	76.97	84.91	79.07	53.85	44.00
Generate	75.15	84.91	74.42	57.69	51.00
SOTA	78.79	88.68	79.07	57.69	55.00
Open Source	72.73	86.79	73.26	42.31	47.00

we find that automatically generated knowledge ("Generate") performs comparably to manually crafted knowledge ("Hand") across most metrics. This suggests that our knowledge acquisition pipeline effectively captures and structures agent experiences, demonstrating that the automated generation of knowledge can ultimately achieve performance comparable to that of manually curated knowledge.

H Related Work

Memory systems have emerged as a fundamental component for enabling continuous learning and adaptability in LLM-based agents [21, 22]. Recent surveys highlight memory’s central role in building brain-inspired, modular agent architectures [4, 64, 65]. However, current memory implementations face critical challenges in efficiently managing large-scale information, maintaining retrieval accuracy, ensuring effective cross-domain knowledge transfer, and preserving readable memory structures [66]. Early memory approaches focused on simple storage mechanisms, but have evolved toward more sophisticated architectures. MEMORYLLM [24] embeds memory within the latent space for post-deployment knowledge injection, but this renders memory opaque. More structured approaches include A-MEM [15], which organizes information in a Zettelkasten-style graph, and AriGraph [25], which integrates semantic knowledge graphs with episodic memory using triplets. Meanwhile, MemGPT [26] implements a hierarchical memory system inspired by virtual memory, while the Unified Mind Model [27] partitions storage based on Global Workspace Theory. Memory systems can be categorized by their architecture and retrieval methods. R2D2 [67] transforms navigation from unknown to known MDPs by analyzing workflow trajectories. HippoRAG [35] mimics hippocampal indexing for efficient retrieval from large episodic stores, while Echo [36] injects temporal cues to enable episodic recall. MemInsight [68] autonomously adds semantic tags to boost retrieval performance, and HiAgent [13] chunks working memory by sub-goals to improve long-horizon task success. These approaches advance beyond basic Retrieval-Augmented Generation (RAG) [34, 69] paradigms. The consolidation of memories—transitioning information from short-term to long-term storage—represents another critical dimension of memory systems. Inspired by neuroscience research on memory consolidation [70], systems like TME [71] and Mobile-Agent-E [72] implement structured frameworks for organizing task-relevant information hierarchically. MARK [30] enhances memory consolidation through specialized agents that analyze stored information, detect patterns, resolve contradictions, and prioritize relevant content based on temporal factors. Despite these advances, a universal memory architecture enabling effective knowledge transfer across domains remains elusive. Second, [73] points toward future directions by proposing an AI-native memory paradigm that serves as an intelligent, persistent memory offload system, autonomously generating context-aware responses and facilitating seamless interaction with external systems. Detailed comparisons are shown in Table 13.

H.1 Agent Knowledge Integration

Effective knowledge integration is essential for enhancing agent planning capabilities and mitigating hallucination issues [74–77]. Research in this area has produced diverse approaches to knowledge representation and utilization. Agent Workflow Memory (AWM) [12] enables automatic induction and reuse of sub-workflows, while KnowAgent [28] augments prompts with action-knowledge bases. Some systems learn parametric world-knowledge models (WKM) [29], while others leverage multi-agent systems for dynamic knowledge adaptation, as in MARK [30]. Generate-on-Graph (GoG) [78] addresses incomplete knowledge graph challenges by allowing agents to generate new factual triples while exploring knowledge graphs, effectively treating LLMs as both agents and knowledge stores. Similarly, Pandora [79] constructs a unified knowledge representation using Python’s Pandas API to align with LLM pre-training, facilitating knowledge transfer across diverse structured reasoning tasks. Studies also demonstrate the value of continually accumulated knowledge [80], state-action trajectories [10], and domain-specific subtask libraries [81]. These approaches significantly improve performance in complex domains, with systems like BioDiscoveryAgent [82] and ChemistX [83] showcasing effectiveness in specialized scientific applications. Case-Based Reasoning (CBR) [41] has emerged as a particularly promising strategy, where agents solve new problems by referencing past experiences, enhancing LLMs with explicit, structured knowledge. This approach enables agents to leverage self-reflection, introspection, and curiosity through goal-driven autonomy mechanisms, thereby creating more sophisticated reasoning capabilities.

H.2 Agentic Workflow Management

Agentic workflow—the high-level planning of agents—has evolved from rigid, pre-defined sequences to adaptive schemes enabling dynamic learning and plan revision [8]. Early systems with “hard-coded” control provided reliability in familiar contexts but proved brittle when facing novel situations [31, 32, 84, 85]. Modern approaches increasingly treat workflow construction as a search or self-

Table 13: Comparison of recent memory systems for LLM-based agents. **Cross Domain**: Whether the memory system supports storing and using knowledge across multiple domains (not limited to one field). **Reusable**: Whether the memory system can be separated from the original agent and used as a general-purpose tool by other agents. **Flexible**: Whether the structure, content format, and retrieval method can be decoupled and flexibly reconfigured or optimized. **Self-eval.**: Whether the system includes pre-analysis of retrieved memory or sub-questions before tool usage or reasoning. **Multi-function**: Whether the agent supports multiple reasoning functions (e.g., plan decomposition, tool use, code generation). **Multi-step**: Whether the memory supports trajectory or plan storage for complex, multi-step tasks.

System	Cross Domain	Reusable	Flexible	Self-eval.	Multi-function	Multi-step
UMM	✓	—	✗	✓	✗	✗
A-MEM	—	✓	✓	✓	✗	—
AriGraph	—	✓	✓	✗	✗	✗
R2D2	✗	✓	✓	✓	✗	✓
HiAgent	✓	✓	✗	✗	✗	✗
Mobile-AE	✓	✓	✗	✓	✗	✓
KnowAgent	✓	✗	✗	✗	✗	✓
AutoPatch	✗	—	✗	✓	✗	✓
AWM	✓	—	✗	✓	✓	✓
ChemAgent	✗	—	✓	✓	✓	✓
AGENT KB	✓	✓	✓	✓	✓	✓

improvement challenge. AFLOW [84] formalizes workflow generation as a Monte-Carlo tree search over code graphs. ReAct [31] synergizes reasoning and acting by interleaving chain-of-thought with tool calls, allowing real-time plan adaptation. Building on this, Reflexion [32] enables agents to learn from verbalized self-critiques of past failures, while Toolformer [33] demonstrates that LLMs can learn to use external tools in an unsupervised manner. Hierarchical systems like MPO [9] and CaPo [86] employ meta-controllers to analyze sub-task performance and adjust high-level strategy. WorkTeam [87] introduces a multi-agent framework comprising supervisor, orchestrator, and filler agents, each with distinct roles that collaboratively enhance workflow construction from natural language instructions. Meanwhile, TME [71] introduces a hierarchical Task Memory Tree structure that tracks execution state using a tree where each node corresponds to a task step, storing relevant input, output, status, and sub-task relationships. TRAD [88] addresses the challenge of selecting and utilizing in-context examples effectively, employing thought retrieval for step-level demonstration selection and aligned decision mechanisms to complement retrieved demonstration steps with preceding or subsequent steps. Synapse [10] introduces state abstraction, trajectory-as-exemplar prompting, and exemplar memory to improve multi-step decision-making and generalization to novel tasks. In memory-augmented embodied agents [39], specialized components like routing agents, task planning agents, and knowledge base agents work together, leveraging in-context learning to avoid explicit model training while using RAG to retrieve context from past interactions.

H.3 Future Directions and Challenges

Despite significant progress, several challenges remain in developing memory and knowledge systems for LLM agents. Shah et al. [14] explore pre-conditions for memory-learning agents, revealing that memory induction and architecture quality significantly impact performance. They observe that transferring memory induced by stronger models can increase success rates, suggesting that effective memory systems have strict capability requirements.

EventWeave [40] addresses incomplete context tracking by identifying both core and supporting events in a dynamic event graph, helping models focus on the most relevant information when generating responses. This approach highlights the need for more sophisticated contextual understanding in long-term agent interactions.

Creating lifelong cognitive systems [66] presents additional challenges in enabling continuous, high-frequency interactions while maintaining incremental learning capabilities. The need to rapidly update with new information while retaining and accurately recalling past experiences remains a significant research frontier.

Key challenges for the future include: (1) integrating parametric and non-parametric memory components for maximum flexibility; (2) developing more efficient memory consolidation mechanisms inspired by neuroscience; (3) creating dynamic knowledge structures that can adapt to changing domains; (4) balancing structure and adaptability in workflow management; and (5) developing more robust evaluation frameworks for memory-augmented systems.

As emphasized by Pink et al. [89], episodic memory represents a critical missing piece for long-term LLM agents, supporting single-shot learning of instance-specific contexts. Approaches like Memento No More [90] point toward solutions where agents internalize knowledge and skills for multiple tasks without relying on expanding prompts, moving beyond the limitations of current memory paradigms.

I Knowledge Base Construction Details

Our AGENT KB implementation incorporates experiences from multiple datasets to ensure comprehensive coverage across different task types and domains. Below, we detail the specific sources used for constructing our AGENT KB. The detailed prompts employed during the generation process are presented in Appendix O.1.

I.1 Experience Construction Process

The experience construction process is a structured methodology aimed at iteratively building, validating, and generalizing task-solving strategies across various domains. This pipeline integrates human expertise with automated learning mechanisms to form a scalable and reusable knowledge base of experiences.

In the first stage, a set of initial experiences is manually designed based on domain-specific knowledge or insights derived from solving the original task. Each experience encapsulates a problem-solving trajectory, including observations, actions, and outcomes. These experiences are then rigorously tested against the original task to assess their effectiveness. Those that successfully resolve or significantly improve performance on the task are retained in an initial experience pool $\mathcal{E}_{\text{initial}} = \{e_1, e_2, \dots, e_n\}$, where each e_i is a tuple (q_i, p_i, x_i) representing a query, plan, and corresponding experience, respectively. Invalid or suboptimal experiences are either refined through iterative testing or discarded.

Once a validated experience pool is established, the next phase involves generalizing these experiences across diverse datasets using prompt-based few-shot learning techniques. A prompt template is constructed to guide the model in generating analogous experiences for new tasks. The prompt includes a query q' , an initial plan p' generated by the model, and a few-shot context containing previously validated examples:

$$\text{Prompt}(q', p') = [q' \| p' \| \{(q_1, p_1, x_1), \dots, (q_k, p_k, x_k)\}]$$

Here, k denotes the number of few-shot examples included in the prompt. Given this input, the language model generates an output x' , which represents the derived experience for the new query:

$$x' = f_\theta(\text{Prompt}(q', p'))$$

where f_θ denotes the parameterized function implemented by the language model. The generated experience x' is then stored alongside its corresponding query and plan in the final experience pool $\mathcal{E}_{\text{final}}$. This expansion enables cross-task generalization, allowing the system to adapt strategies from one domain to another.

To enable efficient reuse of these experiences in future tasks, they are organized into a knowledge base service accessible via an API. All collected experiences are indexed using vector representations $\phi(e) \in \mathbb{R}^d$, obtained through an embedding model $\phi(\cdot)$. When a new query q_{new} is received, it is similarly embedded as $\phi(q_{\text{new}})$, and the system retrieves the most relevant experiences by computing cosine similarity:

$$\text{sim}(q_{\text{new}}, e) = \frac{\phi(q_{\text{new}}) \cdot \phi(e)}{\|\phi(q_{\text{new}})\| \|\phi(e)\|}$$

Top- k experiences with the highest similarity scores are returned via an API endpoint, enabling downstream agents or models to incorporate prior knowledge into their planning or decision-making process. This retrieval-augmented framework ensures that the system remains adaptive and responsive to novel but related tasks, forming a closed-loop experience learning cycle.

I.2 Hand-crafted Experience Process

The procedure of hand-crafted experience is described as follows:

- **Step 1: Team Setup and Objective Definition**

Three computer science students familiar with the GAIA benchmark and agent reasoning workflows were recruited to collaboratively design high-quality prompts. The main objective was to transform successful agent reasoning paths into structured, human-readable instructions that captured essential steps, tools, and decision rules.

- **Step 2: Review of Historical Logs**

Each student was assigned a subset of GAIA benchmark tasks (Level 1, 2, 3). They thoroughly examined the corresponding smolagent logs, focusing on:

- Tasks where the agent reached the correct answer.
- Action sequences that were logically sound and tool-use efficient.
- Common patterns across multiple tasks.

After that, they also analyzed the logs of the failed questions, trying to fix the wrong answers by hand with the successful experience.

- **Step 3: Prompt Authoring and Standardization**

The team synthesized these findings into general reasoning workflows—abstract sequences that could be reused.

Each reasoning pattern was rewritten into a natural language instructional prompt. Prompts were standardized to use consistent sentence structures, imperative voice, and tool-neutral references.

- **Hand Crafted Example Experience:**

Search for the 2015 paper “Pie Menus or Linear Menus, Which Is Better?” on a scholarly database (e.g., Google Scholar or IEEE Xplore) and note the authors in “First M. Last” format. For each author, look up their publication history on DBLP or Google Scholar and list all their papers with publication years. Determine which author has works published before 2015, and collect that author’s prior publications. Sort the author’s earlier papers by year and identify the very first one. Verify the title of that earliest paper against the database entry to ensure accuracy.

- **Step 4: Effectiveness Testing and Selection**

To evaluate quality, each handcrafted experience was tested via few-shot prompting on similar GAIA tasks.

The top 80 prompts with the best performance were selected as the canonical set.

- **Step 5: Generalization to Other Benchmarks**

Using these 80 high-quality examples, we applied few-shot prompting to generate experience instructions for other reasoning benchmarks.

I.3 Experience Source Overview

Our AGENT KB is constructed from a diverse set of benchmark datasets spanning code reasoning, web navigation, multi-hop retrieval, and human-level evaluation tasks. Each dataset contributes structured experience entries that reflect distinct problem-solving patterns and domain characteristics.

Table 14 summarizes the data sources, their original task counts, and the number of resulting experience entries after processing:

Table 14: Overview of datasets used to construct the experience knowledge base.

Dataset	Domain	Original Tasks	Generated Experience
BrowseComp	Web navigation, info retrieval	1,266	1,266
MultiHopRAG	Multi-hop reasoning	2,556	2,556
HLE	Human-level evaluation	3,000	2,000
WebWalkerQA	Open-domain QA	680	680
RepoClassBench	Code understanding, classification	100 (Python subset)	1,000
SWE-Gym-Raw	Code generation, bug fixing	100	1,000
RepoEval	Code completion, repository tasks	100	1,000

I.3.1 GAIA Experience Details

BrowseComp. We processed all 1,266 tasks from the BrowseComp benchmark (https://huggingface.co/datasets/smolagents/browse_comp), creating one experience entry per task. These experiences capture web browsing, information retrieval, and multimodal reasoning patterns.

MultiHopRAG. We incorporated all 2,556 tasks from the MultiHopRAG dataset (<https://github.com/yixuantt/MultiHop-RAG/tree/main/dataset>), with each task contributing one experience entry. MultiHopRAG experiences focus on multi-hop reasoning and retrieval-augmented generation scenarios.

HLE. From the HLE benchmark’s 3,000 tasks (<https://huggingface.co/datasets/cais/hle>), we selected the text-based subset, creating one experience entry per task. We excluded non-textual tasks to maintain consistency in knowledge representation. These experiences cover human-level evaluation scenarios across diverse domains.

WebWalkerQA. We integrated 680 tasks from WebWalkerQA (<https://huggingface.co/datasets/callanwu/WebWalkerQA>), with each task contributing one experience entry. These experiences capture web navigation and question-answering patterns in open-domain contexts.

The resulting knowledge base provides a rich, diverse foundation of experiences spanning programming, reasoning, retrieval, navigation, and multimodal tasks. This diversity is essential for enabling effective knowledge transfer across domains and agent frameworks.

I.3.2 SWE-bench Experience Details

RepoClassBench. We utilized the RepoClassBench dataset (<https://github.com/microsoft/repoClassBench>), selecting 100 representative cases from Python repositories that align with those in SWE-bench. For each case, we generated 10 distinct experiences capturing different solution approaches, resulting in 1,000 structured knowledge entries. These experiences focus on repository classification and code understanding tasks.

SWE-Gym-Raw. We incorporated the SWE-Gym-Raw dataset (<https://huggingface.co/datasets/SWE-Gym/SWE-Gym-Raw>), from which we selected 100 diverse problem instances. Following a methodology similar to RepoClassBench, we generated 10 distinct experiences per instance, resulting in a total of 1,000 knowledge entries. These experiences primarily focus on code generation and bug-fixing scenarios within Python-based repositories.

RepoEval. From the RepoEval dataset (<https://github.com/microsoft/CodeT/tree/main/RepoCoder/datasets>), we selected 100 cases and generated 10 experiences per case, creating an additional 1,000 knowledge entries. RepoEval experiences focus on code completion and repository-level programming tasks in Python.

J Implementation Details

J.1 Evaluation Protocol

Success rate is defined uniformly across both benchmarks as the percentage of tasks completed correctly in a single attempt. For SWE-bench, this corresponds to code repair tasks that successfully resolve the issue and pass all test cases. For GAIA, success requires generating the correct response according to the benchmark’s evaluation criteria. This standardized metric enables fair comparison across tasks and domains while providing a clear assessment of practical utility.

Especially, for the GAIA benchmark, we evaluate on the validation set and report results across three difficulty levels that reflect progressively more complex reasoning and tool usage requirements. Level 1 (53 tasks) includes basic tasks requiring simple reasoning or straightforward information retrieval. Level 2 (86 tasks) involves intermediate complexity, typically requiring multi-step reasoning or the use of external tools. Level 3 (26 tasks) consists of advanced tasks demanding sophisticated reasoning, multiple tool interactions, or specialized domain knowledge. The overall performance is measured using an unweighted average across all 165 tasks, ensuring a balanced evaluation without favoring any specific difficulty level. In our evaluation, we adopt both Pass@1 and Pass@3 metrics to assess different configurations. Pass@1 measures whether the model’s first attempt yields a correct solution, and is used as the primary metric for the `+AGENT KB` and `+AGENT KB ✓` configurations. Pass@3, on the other hand, evaluates whether any of the three independently generated solutions is correct, and is applied to the `+AGENT KB ✓♥` configuration.

For the SWE-bench benchmark, we evaluate the performance on the test set, which comprises 300 issue-pull request pairs from 11 popular Python repositories. The evaluation is conducted through unit test verification, with post-PR behavior serving as the reference solution. Similar to the evaluation protocol used for the GAIA benchmark, we employ both Pass@1 and Pass@3 metrics to assess different configurations. The Pass@1 metric measures whether the model’s first attempt results in a correct solution and is used as the primary metric for the `+AGENT KB` and `+AGENT KB ✓` configurations. Conversely, the Pass@3 metric evaluates whether any of the three independently generated solutions is correct and is applied to the `+AGENT KB ✓♥` configuration.

J.2 Implementation Details

In GAIA and SWE-bench experimental evaluation, we report Pass@1 results for both the `+AGENT KB` and `+AGENT KB ✓` configurations. These metrics reflect the accuracy of the first successful solution generated by the student agent, either its initial attempt or after receiving feedback in the second round. Notably, in evaluating the final results, we treat solutions as correct if they are semantically equivalent, even if they differ in formatting. This evaluation approach ensures a fairer assessment of model performance by avoiding incorrect judgments caused by formatting discrepancies or synonymous output variations. The `+AGENT KB ✓♥` configuration introduces further refinements and is evaluated using the Pass@3 metric. This choice aligns with standard practice in recent public evaluations, as most entries on the current GAIA leaderboard report Pass@3 rather than Pass@1 results. To ensure a fair comparison with these established results, we adopt the widely used Pass@3 metric. Pass@3 measures whether any of the three independently generated reasoning attempts contains a correct solution. Importantly, this metric does not rely on aggregation via voting; instead, it reflects the extent to which correct solutions are covered across multiple diverse reasoning paths. By adopting Pass@3, we ensure better comparability with existing methods while gaining insights into the system’s ability to explore alternative problem-solving strategies, thereby increasing the likelihood of generating at least one valid solution among several attempts. By comparing Pass@1 and Pass@3 results across these configurations, we are able to assess the individual contributions of supervision signals, multi-attempt exploration, and knowledge refinement mechanisms within our framework.

J.3 Smolagents Details

The smolagents framework utilized in our experiments is based on an enhanced version of the original smolagents platform, which provides a solid foundation for agent-based reasoning and problem-solving. This improved version integrates multi-source search tools and a multimodal processing toolbox, significantly enhancing the agent’s capacity to retrieve, interpret, and utilize external

knowledge from heterogeneous sources. These extensions facilitate more effective handling of visual content, structured data, and mixed-format inputs, thereby broadening the agent’s applicability to real-world, knowledge-intensive tasks.

In the smolagents + `+AGENT KB` configuration, we further extend the framework by incorporating our proposed AGENT KB system, while retaining its core search and reasoning mechanisms. This integration demonstrates the flexibility and compatibility of AGENT KB with existing agent architectures, underscoring its potential as a modular enhancement across diverse agent frameworks. Experimental results indicate that the synergy between smolagents’ robust reasoning capabilities and AGENT KB’s knowledge-sharing infrastructure leads to performance gains that surpass the individual capabilities of either system alone, highlighting the value of integrating structured knowledge management into agent-based problem solving.

J.4 OpenHands Details

We develop our AGENT KB methods based on the OpenHands framework, which is a platform for the development of powerful and flexible AI agents that interact with the world in similar ways to those of a human developer. Specifically, the `+AGENT KB` configuration adopts a two-stage teacher-student framework. In the first stage, the student agent attempts the task independently. The teacher then reviews the solution and retrieves relevant prior cases from AGENT KB. Then the teacher provides feedback to the student without access to ground-truth correctness (unsupervised). The student subsequently revises their solution based on this feedback.

J.5 Hyperparameters

All experiments used consistent hyperparameters unless otherwise specified for ablation studies. We set the temperature to 1.0 to maintain a balance between creativity and coherence in model outputs. The top-k sampling parameter was fixed at 0.1, allowing for focused yet diverse generation. For retrieval operations, we used a default top-k value of 3, retrieving the three most relevant knowledge entries for each query. Maximum tokens per response were limited to 4096, providing sufficient space for comprehensive reasoning while maintaining computational efficiency. Knowledge base entries were standardized to 512 tokens, enabling rich representation while controlling storage requirements.

These hyperparameter selections were determined through preliminary experiments to optimize the balance between performance and computational efficiency. We maintained these settings consistently across experimental configurations to ensure fair comparison, with the exception of specific ablation studies designed to evaluate parameter sensitivity.

In our AGENT KB (Knowledge Base) implementation, we evaluate three distinct configurations to systematically assess different components of our knowledge transfer and self-improvement framework. These configurations are designed to explore how supervision, consensus-based reasoning, and formatting consistency influence the overall performance and robustness of agent-based learning.

K Agent Framework

K.1 Multi-Agent Collaboration System of smolagents

In our *smolagents* framework, we extend the *smolagents* multi-agent collaboration system, which comprises two core agents: a management agent (**CodeAgent**) and a retrieval agent (**SearchAgent**). The agent tool capabilities have been comprehensively upgraded to support complex tasks. The system adopts a master-slave collaboration architecture, where **CodeAgent** acts as the main agent overseeing global planning and coordination, while **SearchAgent** is responsible for specialized information retrieval. Communication between the agents is established through standardized interfaces, ensuring both modularity and extensibility.

K.1.1 Agent Tool Capabilities

CodeAgent. The CodeAgent incorporates a set of specialized tools to support task planning, multi-tool orchestration, and code analysis. These include the `VisualInspectorTool`, which is responsible for processing image files; the `AudioInspectorTool`, designed for audio file analysis; and the `TextInspectorTool`, which handles text-based content.

SearchAgent. The SearchAgent, which is primarily responsible for information retrieval tasks, includes the `CrawlerReadTool` for accessing and extracting web content, the `CrawlerArchiveSearchTool` for conducting searches within archived web data, and the `TextInspectorTool`, which supports textual analysis as part of the retrieval and processing pipeline.

The integration of these tools enhances the functional capabilities of individual agents, enabling them to effectively process diverse input modalities and information sources. The following part provides a comprehensive overview of the various tools incorporated into the agents.

`TextInspectorTool` reads non-image files (HTML, PDF, DOCX, XML/CSV, etc.) and converts the raw bytes into clean Markdown, and (optionally) lets the LLM answer a follow-up question about that content; it contains specialized libraries to handle specific formats (such as `openpyxl` for Excel, `python-pptx` for PPT) and automatically trims huge documents to a safe token limit.

`AudioInspectorTool` is a thin wrapper around OpenAI Whisper: it validates that the path ends in `.mp3/.m4a/.wav`, transcribes the audio, then either returns the transcript or feeds it (plus a user question) to an LLM for a structured three-part response.

`VisualInspectorTool` handles `.jpg/.png/.gif/.bmp` (or HTTP image URLs), and asks a multimodal model—Qwen-VL first, GPT-4o as a fallback—to generate a rich description or answer an image-based query, with automatic resizing if the payload is too large.

`WikiSearchTool` is a lightweight Wikipedia client: given a query, it hits the MediaWiki API, follows redirects, and returns title, canonical URL, and intro extract, optionally rewriting the query through a `SearchReflector` for better recall.

`SearchTool` unifies several web engines behind one interface; depending on the `search_type` flag it can call ‘wiki’, ‘google’, ‘baidu’, ‘bing’ and ‘duckduckgo’, supports year filtering, returns a numbered snippet list, and tracks visit history to warn about recently opened URLs.

`CrawlerReadTool` is an async “get-page” utility that downloads full HTML (or other text formats) directly through a `SimpleCrawler`, stripping boilerplate so downstream tools can inspect clean text without a headless browser.

`CrawlerArchiveSearchTool` works similarly but targets archived versions of pages (Wayback, Common Crawl, etc.), enabling the agent to fetch content that has been taken down or changed since publication, which is vital for timeline-sensitive questions or provenance checks.

Task Distribution Workflow. The *smolagents* framework includes a structured task distribution workflow that begins with the `CodeAgent`, which receives and analyzes user-defined tasks. During the analysis phase, the `CodeAgent` identifies subtasks requiring external information retrieval and delegates these to the `SearchAgent`. The `SearchAgent` then performs targeted searches to gather relevant data from appropriate sources. Once retrieved, the information is returned to the `CodeAgent`, which integrates the results into its execution flow to proceed with task completion.

Communication Protocol. To enhance performance and efficiency, smolagents supports several collaboration optimization techniques. These include task priority management to ensure timely execution of critical operations, as well as optimized resource scheduling to improve utilization of available computational resources. The system allows for the concurrent execution of independent tasks to reduce overall processing time. Moreover, a result caching mechanism is integrated to store and reuse previously obtained results, thereby minimizing redundant computations and improving response efficiency for recurring queries.

This structured collaboration mechanism ensures clarity and efficiency in agent interactions, supporting complex, multi-step reasoning while remaining modular and scalable.

L Execution Process Example

Initially, the AGENT KB independently processed a query, leveraging its internal retrieval mechanisms to summarize relevant information. This initial processing led to the generation of a preliminary plan. As part of this plan, the AGENT KB initiated a search for the two most commonly associated chemicals. However, the search successfully retrieved the European Community (EC) number for only one of these chemicals, as the other was not an enzyme. Based on this incomplete retrieval, the AGENT KB incorrectly concluded that the true answer was solely confined to entry "3.1.3.1," leading to two sequential and erroneous returns.

Subsequently, the Teacher Agent provided a more explicit and refined planning directive, offering a clearer understanding of the task trajectory and expected outcome with fewer steps. It clarified that the question was specifically about enzymes, since only enzymes have EC numbers. As a result, the AGENT KB would no longer search for non-enzyme chemicals in the first place, and it would avoid mistakenly returning a "no correct answer" response in the following steps. This guidance helped the AGENT KB overcome the limitations of its initial independent processing and enabled it to retain and utilize the corrected context.

Agent KB – User Question

What are the EC numbers of the two most commonly used chemicals for the virus testing method in the paper about SPF MV and SPCSV in the Pearl Of Africa from 2016? Return the semicolon-separated numbers in the order of the alphabetized chemicals.

Agent KB – Reason the question to get suggestion

EC numbers; virus testing method; SPF MV; SPCSV; Pearl Of Africa; 2016 paper; commonly used chemicals; chemical alphabetization

Agent KB – Retrieve from AgentKB

Similar task: During studies of the phage genome database, the researchers identified an operon.....

Suggestions: Explicitly map experimental design to logical variables—use data frames or grid tables for clarity.....

Agent KB – Student's Initial Plan after refinement

1. Identify the specific 2016 paper on SPF MV and SPCSV in Uganda.
2. Retrieve the full text using inspect_file_as_text or search agent.
3. Check the methodology section for the virus testing method (e.g., ELISA).
4. Extract chemicals used, focusing on enzymes.
5. Determine the two most commonly used ones.
6. Alphabetize the chemical names.
7. Look up their EC numbers using BRENDA, ExplorEnz, etc.
8. Verify EC number correctness.
9. Return EC numbers in the correct order, separated by a semicolon.

Agent KB – Model Prediction

3.1.3.1;3.1.3.1

Agent KB - Retrieve from AgentKB

Similar task: Which of the following peptides has been found to bind or is highly likely to bind an integrin receptor in an in vitro assay.....

Suggestions: Always decompose the question to clarify the biological context and functional motifs encountered in sequence-based tasks.....

Agent KB - Teacher Agent's suggestion

1. Search for the paper using keywords like "SPFMV", "SPCSV", "Pearl of Africa", and "2016".
2. Identify the virus detection method (typically ELISA) and the **enzymes** used in the labeling step.....

Agent KB - Student's New Plan

1. Use compact keywords to Search for the 2016 paper on SPFMV and SPCSV in Pearl of Africa.
2. Review the method section to find the virus detection technique and the two **enzymes**.
3. Look up EC numbers from BRENDA or ExPASy.
4. Alphabetize enzyme names.
5. Return EC numbers joined by a semicolon.

Agent KB - Model Thought

I suspect the two most common enzymes for virus detection are alkaline phosphatase and horseradish peroxidase. Their EC numbers are:

- Alkaline phosphatase: 3.1.3.1
 - Horseradish peroxidase: 1.11.1.7
- Alphabetically, alkaline phosphatase comes first.
So I will return: 3.1.3.1;1.11.1.7

Agent KB - Final Output (Code)

```
answer = "3.1.3.1;1.11.1.7"  
final_answer(answer)
```

Agent KB - Observation

Last output from code snippet:
3.1.3.1;1.11.1.7

Table 15: AGENT KB Workflow Example of SWE-bench

Agent KB - Problem Description
<p>In the project that automatically generates API documentation for Python projects, an extra backslash is inserted before underscores—for example, hello_ is rendered as hello_</p>
Agent KB - Retrieve from AgentKB
<p>Check the string processing part in the relevant functions to ensure that escape is only carried out when necessary. For example, whether all parameters ending with _ need to be escaped, or whether different handling methods are required in certain specific contexts (such as attribute names, parameter names, etc.).</p> <p>When modifying conditions, not only the original conditions should be taken into account, but also factors such as configuration and context should be combined to ensure the accuracy of the logic.</p>
Agent KB - Problem Description
<p>Disabling evaluation globally with with evaluate(False) interferes with sympify's string-parsing logic, preventing some integer expressions from being instantiated as integer objects.</p>
Agent KB - Retrieve from AgentKB
<p>When adding or modifying a conditional check (such as for 'evaluate' or imaginary coordinates), ensure the logic does not inadvertently skip important validation for invalid inputs (such as actual imaginary numbers), and only disables overly strict checks for valid real inputs. This is critical to maintain mathematical correctness while fixing the bug. (Most important)</p> <p>When changing the logic in constructors (like Point/Point2D), verify that the minimal change solves the immediate bug, does not introduce new regressions, and does not allow forbidden cases (e.g., actual imaginary coordinates)</p>

M Additional Discussion

Our experimental results demonstrate the effectiveness of AGENT KB in enhancing agent performance across a range of tasks and model architectures. The significant improvements observed, particularly for medium-difficulty tasks, suggest that shared memory infrastructures can play a crucial role in advancing agent capabilities.

Several aspects of our findings warrant further discussion:

1. **Scalability:** As the knowledge base grows, efficient retrieval becomes increasingly challenging. Future work should explore more sophisticated indexing and retrieval mechanisms to maintain performance at scale.
2. **Diversity:** The diversity of knowledge sources is critical for robust performance. Expanding AGENT KB to include a wider range of tasks and domains could further enhance its generalization capabilities.
3. **Quality control:** As community contributions grow, maintaining the quality and reliability of the knowledge base becomes more challenging. Developing automated verification and validation mechanisms will be essential.
4. **Adaptation:** Our current approach relies primarily on retrieval and reuse of existing knowledge. Developing more sophisticated adaptation mechanisms could enhance the system's ability to handle novel tasks and contexts.

N Agent KB Data Examples

This section provides concrete examples of how AGENT KB processes and stores different types of agent experiences. We demonstrate three key components: SWE-bench workflow examples showing problem-solution pairs, raw execution logs transformed into structured experiences, and complex multi-constraint query processing.

N.1 SWE-bench Workflow Examples

The following examples illustrate how AGENT KB stores and retrieves domain-specific knowledge for software engineering tasks from the SWE-bench dataset. Each example shows a problem description paired with relevant guidance retrieved from AGENT KB, demonstrating the system's ability to provide contextual assistance for code debugging and modification tasks.

Table 16: AGENT KB Workflow Example of SWE-bench

Agent KB - Problem Description In the project that automatically generates API documentation for Python projects, an extra backslash is inserted before underscores—for example, hello_ is rendered as hello_. Agent KB - Retrieve from AgentKB Check the string processing part in the relevant functions to ensure that escape is only carried out when necessary. For example, whether all parameters ending with _ need to be escaped, or whether different handling methods are required in certain specific contexts (such as attribute names, parameter names, etc.). When modifying conditions, not only the original conditions should be taken into account, but also factors such as configuration and context should be combined to ensure the accuracy of the logic. Agent KB - Problem Description Disabling evaluation globally with with evaluate(False) interferes with sympify's string-parsing logic, preventing some integer expressions from being instantiated as integer objects. Agent KB - Retrieve from AgentKB When adding or modifying a conditional check (such as for 'evaluate' or imaginary coordinates), ensure the logic does not inadvertently skip important validation for invalid inputs (such as actual imaginary numbers), and only disables overly strict checks for valid real inputs. This is critical to maintain mathematical correctness while fixing the bug. (Most important) When changing the logic in constructors (like Point/Point2D), verify that the minimal change solves the immediate bug, does not introduce new regressions, and does not allow forbidden cases (e.g., actual imaginary coordinates)
--

These examples demonstrate AGENT KB's ability to provide targeted guidance for common software engineering challenges. The first example addresses API documentation generation issues with string escaping, while the second focuses on symbolic mathematics library debugging. Notice how the retrieved knowledge provides specific, actionable advice rather than generic troubleshooting steps.

N.2 Raw Log to Experience Generation

This subsection demonstrates the complete pipeline for transforming raw agent execution logs into structured knowledge that can be stored in AGENT KB. This process is crucial for the system's learning capability, allowing successful problem-solving strategies to be captured and reused.

N.2.1 Raw Log Example

The following demonstrates how agent execution logs are processed and transformed into structured experiences for AGENT KB. This particular example shows a bioinformatics task involving protein structure analysis, where the agent had to adapt its approach when encountering unexpected file formats.

```
1 {
2     "agent_name": "gpt-4.1",
3     "question": "Using the Biopython library in Python, parse the PDB file of the
4         protein identified by the PDB ID 5wb7 from the RCSB Protein Data Bank. Calculate
5         the distance between the first and second atoms as they are listed in the PDB
6         file. Report the answer in Angstroms, rounded to the nearest picometer.",
7     "prediction": "1.46",
8     "true_answer": "1.456",
9     "intermediate_steps": [
10        {
11            "task": "You have one question to answer...",
12            "step_type": "task"
13        },
14        {
15            "facts": "Here are the facts that I know so far...",
16            "plan": "Here is the plan of action that I will follow...",
17            "step_type": "planning"
18        },
19        {
20            "tool_calls": [{"id": "call_1", "type": "function", "function": {"name": "python_interpreter", "arguments": "..."}}],
21            "error": {"type": "AgentExecutionError", "message": "Code execution failed..."},
22            "step_type": "action"
23        }
24    ]
25 }
```

Listing 1: Raw Agent Execution Log

N.2.2 Key Insights

These examples collectively demonstrate several important aspects of AGENT KB's design and functionality:

1. **Domain Adaptation:** The system successfully captures domain-specific knowledge across different fields (software engineering, bioinformatics, biographical research), showing its general applicability.
2. **Error Recovery:** Raw logs show how agents adapt when initial approaches fail, and these adaptation strategies are preserved as valuable experiences for future use.
3. **Precision Management:** The system learns specific formatting and precision requirements, crucial for tasks requiring exact numerical outputs.
4. **Multi-Agent Coordination:** Complex queries demonstrate how different agent roles (general reasoning vs. specialized search) can be coordinated with distinct but complementary planning strategies.
5. **Source Validation:** The emphasis on authoritative sources and cross-validation shows the system's commitment to reliability and accuracy in information gathering.

These examples illustrate how AGENT KB transforms individual agent experiences into a shared knowledge resource that enhances the performance of the entire agent ecosystem.

N.2.3 Generated Agent Experience

From the raw log, the following agent experience is extracted:

Agent Experience

When calculating distance, extract only the first two ATOM coordinates directly without skipping any lines. If the task asks for the distance in Ångströms, rounded to the nearest picometer, keep the original Å value with three decimal places precision, without converting back and forth between Å and pm. Output the distance directly in Ångströms, keeping the computed value with three decimals (e.g., x.xxx), do not round it to fewer decimal places. The output is just the number without any units or symbols, such as x.xxx, not x.xxx Ångströms.

N.3 Complex Query Processing Example

N.3.1 Question-Answer Format

Complex Query Example

Question: There's this popular figure with multiple Grammy awards in the entertainment industry, who put out their first album before 1969, had substance dependence, and was dismissed from school before they turned 20. Their first life partner died in 1997. They became uniformed personnel at some point in their lives. What's the name of the hospital they died?

True Answer: St. John's Health Center

N.3.2 Generated JSON Structure

The following JSON structure shows how AGENT KB organizes planning strategies and experiences for complex queries. The structure includes separate planning approaches for different agent types (general agents vs. specialized search agents) and captures both successful strategies and lessons learned from the task execution. This structured approach enables systematic knowledge transfer and strategy refinement.

```
1 {
2   "question": "There's this popular figure with multiple Grammy awards...",
3   "true_answer": "St. John's Health Center",
4   "agent_planning": "1. Parse the question to extract all key constraints: multiple
      Grammy awards, first album before 1969, substance dependence, dismissed from
      school before age 20, first life partner died in 1997, served as uniformed
      personnel, determine place/hospital of death.\n2. Conceptual plan:\n- Identify
      the possible entertainers matching all constraints.\n- For each candidate:\n  a)
      Verify the timeline for first album release (before 1969)\n  b) Check Grammy
      history\n  c) Search biographical records for substance abuse and educational
      background\n  d) Confirm information about life partner's death year and
      uniformed service\n  e) Pinpoint the date and location/hospital of death of the
      matched figure.",
5   "search_agent_planning": "1. Receive precise person identifier from Code Agent or
      use biographical clues to triangulate the subject.\n2. Formulate search queries
      for identification and specific hospital information.\n3. Prioritize official
      biographical sources, reputable news outlets, Grammy records.\n4. Cross-check
      critical data points to validate subject match.\n5. Extract facts about location
      and hospital of death from obituaries.",
6   "agent_experience": [
7     "Break down multifaceted questions into smaller constraint checks",
8     "Explicitly log and verify biographical constraints with multi-source
       confirmation",
9     "Select high-reliability sources for biographical and award data",
10    "Delegate to Search Agent early with specific sub-queries",
11    "Validate final answers by chaining all found facts back to original constraints"
12  ],
13  "search_agent_experience": [
14    "Decompose complex queries into sequential search refinements",
```

```
15     "Craft highly specific queries for ambiguous identifiers",
16     "Favor authoritative sources over entertainment/tabloid content",
17     "Cross-validate information from multiple independent sources",
18     "Format results with direct attribution and clear source references"
19 ]
20 }
```

Listing 2: Generated Agent Planning and Experience JSON

O Collections of Used Prompts

O.1 Prompt Design for AGENT KB Database Generation

O.1.1 General Tasks

Agent KB Generation Prompt

You will act as an advanced AI evaluation system tasked with analyzing a complex problem handled by an agent. Your analysis will extract valuable insights from this process. Follow these instructions carefully:

1. I will provide a question and its correct answer (true_answer).
2. First, simulate the agent's planning process in detail. Describe how it would:

- Break down the problem into logical components
- Determine which tools to use (code execution, data processing, API calls)
- Decide when to delegate to the Search Agent
- Plan data transformations and analysis steps
- Structure the final solution

Include specific reasoning steps, potential code snippets considered, and decision points. Only include content to the agent plan, without any other description.

4. Next, based on the question and your simulated planning processes, create a realistic error scenario. Describe:

- Where and how the agents might fail
- Incorrect assumptions they might make
- Data misinterpretations or code errors
- Logical flaws in their approach

5. Finally, provide actionable experience guidelines:

- Specific principles to improve problem-solving, tool selection, verification, and integration of search results

The behavioral guidelines should be generalizable principles that would help the agents perform better on similar tasks, without directly revealing the specific answer to the question I provided.

Output your complete analysis in the following JSON format with no additional text:

```
{  
  "question": "<question I provide>",  
  "true_answer": "<correct answer I provide>",  
  "agent_plan": "<your detailed Code Agent plan simulation>",  
  "agent_experience": "<your actionable Code Agent guidelines>",  
}
```

Here is an example:

```
{  
  "question": "<question from hand-crafted experience pool>",  
  "true_answer": "<correct answer>",  
  "agent_plan": "<Real Code Agent plan>",  
  "agent_experience": "<Hand-crafted agent experience>",  
}
```

0.1.2 GAIA

Agent KB Generation Prompt

You will act as an advanced AI evaluation system tasked with analyzing a complex problem handled by a Code Agent with an embedded Search Agent. Your analysis will extract valuable insights from this process. Follow these instructions carefully:

1. I will provide a question and its correct answer (true_answer).
2. First, simulate the Code Agent's planning process in detail. Describe how it would:

- Break down the problem into logical components
- Determine which tools to use (code execution, data processing, API calls)
- Decide when to delegate to the Search Agent
- Plan data transformations and analysis steps
- Structure the final solution

Include specific reasoning steps, potential code snippets considered, and decision points. Only include content to the agent plan, without any other description.

3. Next, simulate the Search Agent's planning process in detail. Describe how it would:

- Parse the search query requirements from the Code Agent
- Formulate effective search queries
- Determine which sources to prioritize
- Extract and validate relevant information
- Process and structure the search results for the Code Agent

Include specific query formulation strategies and information filtering approaches. Only include content to search agent plan, without any other description.

4. Based on the question and your simulated planning processes, create a realistic error scenario. Describe:

- Where and how the agents might fail
- Incorrect assumptions they might make
- Data misinterpretations or code errors
- Logical flaws in their approach

5. Finally, provide two sets of actionable experience guidelines:

- For the Code Agent: Specific principles to improve problem-solving, tool selection, verification, and integration of search results
- For the Search Agent: Specific principles to enhance query formulation, source evaluation, information extraction, and result formatting

The behavioral guidelines should be generalizable principles that would help the agents perform better on similar tasks, without directly revealing the specific answer to the question I provided.

Important: If the question does not require the search agent to solve, leave "search_agent_plan" and "search_agent_experience" empty in your response.

Output your complete analysis in the following JSON format with no additional text:

```
{  
  "question": "<question I provide>",  
  "true_answer": "<correct answer I provide>,"
```

```

"agent_plan": "<your detailed Code Agent plan simulation>",
"search_agent_plan": "<your detailed Search Agent plan simulation>",
"agent_experience": "<your actionable Code Agent guidelines>",
"search_agent_experience": "<your actionable Search Agent guidelines>"
}

```

Here is an example:

```

{
"question": "<question from hand-crafted experience pool>",
"true_answer": "<correct answer>",
"agent_plan": "<Real Code Agent plan>",
"search_agent_plan": "<Real Search Agent plan>",
"agent_experience": "<Hand-crafted agent experience>",
"search_agent_experience": "<Hand-crafted search agent experience>"
}

```

0.1.3 SWE-bench

Agent KB Generation Prompt

You are an advanced code repair analysis system tasked with constructing structured experiences for Agent KB from SWE-bench tasks. Given a natural language problem description, a model-generated fix, and supporting repair hints, follow the steps below to extract reusable knowledge entries. Your output should conform strictly to JSON formatting and follow the key structure outlined in each step.

1. Code Reconstruction:

Given a detailed natural language description of a Python class or function, generate its correct implementation. Ensure it is complete and syntactically valid.

Output key: "code"

2. Error Analysis and Repair Principles:

You are given two versions of code: one with errors and one corrected. Analyze the differences and identify key problems in the faulty version. Based on this comparison, produce a list of 10 code repair precautions. These should be generalizable principles addressing common issues (e.g., indentation, type conversion, exception handling, logic errors). Avoid using titles; just output the explanations.

Output key: "hints" (as a list of 10 strings)

3. Hint Classification:

Each natural language hint is used to prompt the LLM to repair the code. Classify each hint into one repair category (e.g., "syntax", "logic", "exception handling"). Also, extract important keywords and write a one-sentence summary of the hint.

Output keys: "category", "keywords", "summary"

4. Repair Type Identification:

Given the original problem description, identify the {K} most relevant categories this code repair case falls under. Select from a pre-defined set of bug types.

Output key: "categories" (as a list of {K} strings)

5. Most Relevant Hints Ranking:

You are given a set of all the hints provided to the model. Analyze the model's generated fix and its reasoning trace. Based on this analysis, identify the {N} most relevant hints. These may be either positively helpful or misleading. Sort them in order of influence on the final patch. Output key: "hints" (as a list of {N} strings)

Important Notes:

- Always respond strictly in JSON format.
- Do not include section titles, markdown formatting, or explanations.
- When code is requested, return only the code inside the JSON key.
- If any step is not applicable (e.g., hint classification not possible), return an empty string or array for that field.

O.2 Prompt Design for AGENT KB Pipeline

O.2.1 GAIA

Agent KB Reason Prompt

Analyze similar tasks and past experiences to generate concise, actionable suggestions for improving the current plan. Based on the patterns identified in relevant tasks and insights from the Agent KB, provide specific recommendations.

Key Requirements:

1. Focus exclusively on technical/behavioral improvements derived from similar task patterns and experience.
2. Provide root-cause solutions and implementation strategies based on past successes.
3. Format output strictly as:
{1. Specific suggestion 1}
{2. Specific suggestion 2}
...

No headings, explanations, or markdown.

You can refer to similar tasks, plans, and corresponding experience to provide your suggestions:

```
{  
  "question": "<Question retrieved from Agent KB>",  
  "agent_plan": "<Retrieved agent plan>",  
  "agent_experience": "<Retrieved agent experience>",  
}  
...
```

Agent KB Refine Prompt

Analyze the execution logs to determine the causes of the agent's incorrect responses. Based on the findings of the log and insights from the provided similar tasks and experience, generate some concise, actionable suggestions that the agent must follow to improve accuracy.

Key Requirements:

1. Focus exclusively on technical/behavioral fixes derived from log patterns and the Agent KB.
2. Provide root-cause resolution (e.g., code logic, data validation, API

handling) as well as generic advice.

3. Format output strictly as:

```
{1. Specific suggestion 1}
{2. Specific suggestion 2}
...
No headings, explanations, or markdown.
```

You can refer to similar tasks and corresponding experience to provide your suggestions:

```
{
  "question": "<Question retrieved from Agent KB>",
  "agent_plan": "<Retrieved agent plan>",
  "agent_experience": "<Retrieved agent experience>",
}
```

...

Execution logs summary:

```
<Log summary>
```

O.2.2 SWE-bench

Agent KB Reason Prompt

Extract key information from user queries to construct efficient search terms for retrieving the most relevant results.

Requirements:

Analyze the user's question to identify core concepts, terminology, and keywords
 Extract contextual information and constraints that may impact search quality
 Break down complex questions into searchable components

Identify the domain, subject matter, and specific needs of the question
 Output format:

```
{<core concepts or topics of the question>}
```

Ensure search terms are specific enough to retrieve relevant information while maintaining sufficient breadth to capture related cases. Combine technical terminology with everyday expressions to optimize search effectiveness.

Agent KB Retrieve Prompt

Given the current bug description, initial patch plan, and model thought process, retrieve the most relevant historical experiences from Agent KB.

Retrieval Priorities:

1. Prefer experiences with similar bug types (e.g., off-by-one errors, null pointer exceptions, wrong return value).
2. Favor patches with successful unit test outcomes and generalizable fix patterns.
3. Include agent plans that show tool usage, exception guards, or correct interface assumptions.

Format each retrieved experience as:

```
{
  "question": "<SWE-bench issue title or commit description>",
  "agent_plan": "<Historical high-level patch or thought process>",
  "agent_experience": "<Failure modes avoided or debug strategies that worked>"}
```

...

Retrieve 3 to 5 relevant entries and return them in the above format for use in downstream reasoning and refinement.

Agent KB **Refine** Prompt

Analyze the execution trace of the model's patch attempt and identify the reasons for its failure. You are given: a natural language description of a code fix problem, the model-generated fix, the model's internal thought process, and the prompts previously provided to guide the model.

Based on this information, identify the most likely cause of the error and determine which hints or prompt components influenced the model's incorrect reasoning. Rank the provided prompts in order of their influence over the model's behavior.

Key Requirements:

1. Focus exclusively on technical root causes, such as incorrect API assumptions, scope misunderstanding, faulty patch structure, or missing validation.
2. Identify which prompt(s) led the model astray, based on reasoning steps or patch behaviors.
3. Output a strictly ranked list of prompts or hints, based on their importance in shaping the erroneous behavior.
4. Justify the ranking based on model thought content and the specific failure observed.

Format strictly as:

```
{  
1. "<Most influential prompt or hint snippet>"  
2. "<Second most influential prompt or hint snippet>"  
...}
```

Do not include headings, explanations, or markdown. Focus only on returning the ranked list with brief justifications inline.