Improving Language Agents through BREW: Bootstrapping expeRientially-learned Environmental knoWledge

Anonymous Author(s)

Affiliation Address email

Abstract

Large Language Model (LLM)-based agents are increasingly applied to tasks requiring structured reasoning, tool use, and environmental adaptation, such as data manipulation, multistep planning, and computer-use automation. However, despite their versatility, current training paradigms for model weight optimization methods, like PPO and GRPO, remain relatively impractical with their high computational overhead for rollout convergence. In addition, the resulting agent policies are difficult to interpret, adapt, or incrementally improve. To address this, we investigate creating and refining structured memory of experiential learning of an agent from its environment as an alternative route to agent optimization. We introduce BREW (Bootstrapping expeRientially-learned Environmental knoWledge), a framework for agent optimization for downstream tasks via KB construction and refinement. In our formulation, we introduce an effective method for partitioning agent memory for more efficient retrieval and refinement. BREW uses task graders and behavior rubrics to learn insights while leveraging state-space search for ensuring robustness from the noise and non-specificity in natural language. Empirical results on real world, domain-grounded benchmarks – OSWorld and τ^2 Bench – show BREW achieves 10-20% improvement in task precision, 10-15% reduction in API/tool calls leading to faster execution time, all while maintaining computational efficiency on par with base models. Unlike prior work where memory is treated as static context, we establish the KB as a modular and controllable substrate for agent optimization – an explicit lever for shaping behavior in a transparent, interpretable, and extensible manner.

23 1 Introduction

2

3

5

6

8

9

10

11

12 13

14

15

16

17

18 19

20

21

22

Large Language Model (LLM) based agents are rapidly being deployed for structured reasoning, tool use, and autonomous interaction in real-world environments [16]. From computer-use and 25 spreadsheet automation to software engineering pipelines, these agents drive tasks such as multi-step planning, data manipulation, and adaptive workflows [21, 13, 32, 2, 19]. For example, a language 27 agent might help automate a multi-step workflow like collecting data from different sources, cleaning 28 or validating it, and then uploading it onto a dedicated server, all while adjusting its plan if the 29 format or structure of the data changes unexpectedly [31, 35, 25, 3]. Yet, despite these successes, top-30 performing agents generally score underwhelmingly on challenging real-world benchmarks—well 31 behind human experts, who routinely exceed 70% success rates [34, 4, 27, 18]. As an example, 32 consider the following scenario:

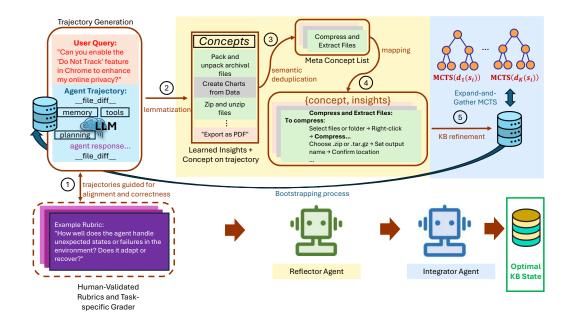


Figure 1: BREW architecture overview using examples from the OSWorld dataset. Step 1 indicates the trajectory generation process with agent alignment to human-validated rubrics and correctness using task-specific grader. Steps 2–4 indicate the Reflector Agent, which learns key concepts and corresponding insights from trajectories. Step 5 indicates the Integrator Agent, which integrates knowledge from the Reflector Agent to bootstrap the KB. We introduce Expand-and-Gather MCTS to further find the best KB configuration as the KB is iteratively refined through reward-guided optimization.

A computer-use agent in an Ubuntu environment tasked with automating software installation across multiple sessions. In its first encounter, it struggles through a 47-step process: opening the wrong package manager, executing redundant dependency checks, and making 23 API calls to complete what could be a 6-step workflow. When presented with a similar installation task in the next session, the agent repeats the same inefficient exploration – *as if encountering the problem for the first time*. A human user, by contrast, would likely have a recollection from internalized memory of the optimal sequence after the first attempt, recognizing the environmental patterns and tool combinations that lead to success.

This scenario illustrates a fundamental limitation of current language agents: despite their impressive capabilities in reasoning and tool use, they lack the ability to accumulate and apply experiential knowledge across task sessions. Each interaction begins from a blank slate, forcing agents to repeatedly explore the same action spaces and rediscover the same solutions [9]. Real-world tasks like long horizon multi-stage automation demand more than just "reactive" [33] tool loops. They require persistent & interpretable learnings from past experiences - what works, what fails and why. To close this gap, recent work has explored learning agent behavior using model weight optimization [23, 22, 24], where agents are trained to maximize success across a wide variety of tool-use episodes. However, while conceptually sound, this suffers from practical limitations. First, it requires expansive exploration over large rollout spaces to converge, especially in domains where tasks are diverse, goals are sparsely defined, and intermediate feedback is noisy or delayed. Second, the resulting policies are often opaque—difficult to interpret, revise, or debug—limiting their real-world deployability. Finally, these policies are tightly coupled to the task distributions they were trained on, making it difficult to adapt or incrementally improve them when downstream requirements shift.

In contrast, others have explored learning of knowledge onto a memory module that remains attached to an agent. These existing memory-augmented agents can be broadly classified into either ones which (i) store only transient trajectory contexts that vanish between episodes like Mem0 [7, 29], or (ii)

embed high-level notes directly in the prompt such as MetaReflection [10] and GEPA[1]. While the latter often do not retain actionable details for future simple tasks, neither of these approach supports modular updates, fine-grained retrieval, or transparent inspection of what the agent "knows." [28].

Leveraging learnings from both camps, we introduce BREW (Bootstrapping experientially-learned environmental knowledge), a framework that incrementally constructs and refines a knowledge base (KB) a structured collection of concept-level documents in natural language, directly from an agent's past interactions. This KB then serves as a persistent memory for the agent to retrieve knowledge in future executions to improve precision and efficiency outcomes. Our key contributions are—

- Novel experience-driven KB construction. We propose a technique for leveraging agent's past interaction trajectories to generate uniquely-partitioned concept-level KB documents. This process is guided by rubrics and task-specific graders which ensures that memories are both semantically aligned with task objectives and human-interpretable.
- State-space search for memory optimization. We formalize the selection and update of KB entries
 as a state search problem and introduce an efficient reward-guided learning scheme, Expand-and Gather Monte Carlo Tree Search (EG-MCTS), that learns to prioritize the most impactful memories
 for robust, multi-step reasoning.
- State-of-the-art results. On domain-grounded benchmarks including OSWorld and τ^2 Bench, BREW achieves significant gains of in the range of 10-20% towards task precision as well as 10-15% fewer steps leading to faster execution, while maintaining memory and compute costs comparable to base LLMs.

72 Preliminary & Related Work

Agent Learning from Demonstrations Recent work has leveraged LLMs to isolate reusable skills through interactive decomposition: one method distills sub-goals from expert trajectories into hierarchical planning and execution policies [11], and another synthesizes executable functional abstractions for advanced mathematical reasoning via program induction [14]. These approaches focus on structured skill extraction from LLM-guided interactions, yet remain reliant on static decomposition or offline synthesis. In contrast, BREW dynamically constructs and refines an experiential memory—learning necessary semantic fragments via rollout-generated insights and structured knowledge-base search (MCTS)—to support long-horizon, memory-augmented planning.

Agentic Memory The concept of providing agents with controllable memory has a rich history. 81 [17]. Memory mechanisms are attracting more and more attention lately [20, 26, 28, 7, 30, 12]. These 82 works focus towards storing relevant context in a structured format like graph or a tree so as to RAG 83 over it. Despite their effectiveness these methods perform well for most cases. However, when the 84 queries are ambiguous, requires multi-hop reasoning and long range comprehension these techniques struggle to perform the tasks [12]. In contrast to prior works BREW uses a state search to explore possible memory states. This allows BREW to select the memory state where the reward during 87 exploration is highest making it more robust to ambiguous queries and long range comprehension. 88 We employ MCTS [8] as a state search algorithm to explore the potential states of the memory by 89 expanding to new and potentially different states of memory based on same interactions. We discuss 90 the state search process more formally in Section 3.3. 91

3 BREW: Architecture

This section describes our proposed **B**ootstrapping expe**R**ientially-learned **E**nvironmental knoWledge model, BREW, which constructs and iteratively refines a KB using trajectory insights guided by human-validated general-purpose agent behavior metrics, task-specific evaluation, and latent insight generation. We introduce a novel decomposition the problem of learning the optimal KB by partitioning memory as local documents associated with semantic concepts, and formulate the KB learning problem as a state space search by proposing Expand-and-Gather Monte Carlo Tree Search (EG-MCTS). Figure 1 provides an architecture overview of BREW, and Algorithm 1 describes pseudocode.

3.1 Trajectory Generation

101

131

Given the training dataset, we generate full-length trajectories, hereby referred to as rollouts, for each query using an LLM-powered agent conditioned on its associated KB. At initialization, the KB is empty, and the LLM is used with a decoding temperature of 0 to ensure deterministic behavior. Further details on training and test splits are in Experiments Section. Each rollout is evaluated using a correctness grader, which assigns a binary label: *success* or *failure* and a qualitative rubric assessment against a set of human-validated general-purpose agent behavior rubrics [6] (Step 1 in Figure 1).

108 3.2 Reflector and Integrator Agents

Reflector Agent: ReflAgent takes as input a rollout with its rubric and correctness labels, and outputs sentence-level insights with mapped concepts:

$$\{concepts, insights\} = \text{ReflAgent}(\{rollout, eval\}).$$
 (1)

Examples of concept-insight pairs appear in Step 2 of Figure 1.

Concept Deduplication: Concept—insight pairs are annotated independently per rollout, often producing overlapping or paraphrased concepts. We address this via semantic clustering (Steps 3–4, Figure 1; Algorithm 1, line 3): contextual embeddings for each concept are generated using an LLM, clustered, and each insight is mapped to its cluster representative. Details appear in Algorithms 2 and 3 in Appendix A.

Integrator Agent: IntegAgent incrementally builds and refines KB documents $\{d(s_i)\}\in\mathcal{D}(s_i)$ during environment interaction. Instead of a centralized memory, the KB is partitioned into local documents, each tied to a meta concept. This design enables (1) efficient, context-specific retrieval; (2) modular updates with minimal interference; and (3) natural alignment with task semantics, as deduplicated meta concepts capture meaningful behavioral abstractions. Unlike prior work assuming flat memory or dialogue histories, this structure is well-suited for long-horizon, procedural tasks where behaviors cluster around discrete skills.

The KB is dynamically populated: concepts central to the dataset receive more updates, shaping memory around frequent behaviors. At each state, for meta concept k, IntegAgent updates its document d_k via

$$d_k(s_{i+1}) \leftarrow \text{IntegAgent}(k, insights_k, d_k(s_i)).$$
 (2)

To reduce LLM variance and improve consistency, we use the Expand-and-Gather MCTS (EG-MCTS) method (Figure 2).

Formally, the KB at state s_i is the union of all concept-localized documents:

$$\mathcal{D}(s_i) = \bigcup_{k \in \mathcal{K}} \{d_k(s_i)\},\tag{3}$$

where K is the set of all meta concepts and $d_k(s_i)$ is the document for concept k at state s_i .

3.3 Expand-and-Gather MCTS for Optimal KB Search

We start by creating a set of meta-concepts after deduplicating concepts extracted by ReflAgent using the first set of trajectory rollouts. We freeze this meta-concept set K, and use it to initialize a KB with an empty document per concept $k \in K$.

We model the problem of finding the optimal KB \mathcal{D}^* as a search problem in the *state* space of all possible KBs \mathcal{D} . To simplify this state search, we model KB \mathcal{D} as a collection of concept level documents. This modeling allows us to break down the larger search space into a collection of simpler document level search problems for each concept k to find the optimal document d_k^* . We then construct the optimal KB \mathcal{D}^* by combining all optimal documents d_k^* for each concept k as follows:

$$\mathcal{D}^* = \bigcup_{\forall k} \{d_k^*\} \tag{4}$$

Notably, even though we are modeling document level search as independent optimization problems, each document in the KB is *not* independent of the others. For example, an agent can retrieve any

Algorithm 1 BREW: Bootstrapping Experientially-learned Environmental Knowledge

Require: Training samples Q_{train} , eval samples Q_{eval} , rubrics, iterations M, candidates per expansion

Ensure: Optimized KB \mathcal{D}^*

```
Initialization
```

```
1: \mathcal{D}_0 \leftarrow \emptyset
 2: \mathcal{B} \leftarrow \text{GENERATEINSIGHTS}(\mathcal{Q}_{\text{train}}, \mathcal{D}_0, \text{rubrics})
 3: \mathcal{K} \leftarrow \text{DEDUPLICATECONCEPTS}(\mathcal{B})
                                                                                                                                               ▷ Initial concept set
 4: for each k \in \mathcal{K} do
             d_k^0 \leftarrow \text{IntegAgent}(k, \mathcal{I}_k, \varnothing)
 5:
             Initialize tree<sub>k</sub> with root node d_k^0
 6:
 7: end for
 8: \mathcal{D}_{\text{current}} \leftarrow \bigcup_{k \in \mathcal{K}} \{d_k^0\}
                                                                                                                                                             EG-MCTS Optimization
 9: for t = 1 to M do
                                                                                                                 ▶ Parallel expansion across concepts
             for each k \in \mathcal{K} do
10:
                    s_k \leftarrow \text{SELECTBESTNODE}(\text{tree}_k)
\mathcal{D}_{\text{best}} \leftarrow \bigcup_{k' \in \mathcal{K}} \{d_{k'}^{\text{best}}\}
\text{EXPANDNODE}(s_k, k, h, \mathcal{D}_{\text{current}}, \mathcal{D}_{\text{best}}, \text{tree}_k)
                                                                                                                                                     11:
                                                                                                                                                12:
13:
14:
             end for
                                                                                                                       ▶ Update current best documents
15:
             for each k \in \mathcal{K} do
                    d_k^{\text{best}} \leftarrow \text{best document in tree}_k
16:
17:
             \mathcal{D}_{\text{current}} \leftarrow \bigcup_{k \in \mathcal{K}} \{d_k^{\text{best}}\}
18:
19: end for
20: return \mathcal{D}_{current}
```

Time Complexity: $O(|\mathcal{Q}_{\text{train}}| \cdot T_{\text{LLM}} + M \cdot |\mathcal{K}| \cdot h \cdot T_{\text{agent}})$

document in the KB during inference and this retrieval making it hard to assess the impact of changing a document in isolation. To solve this we propose Expand-and-Gather MCTS (EG-MCTS), which 143 enables searching these disjoint state spaces concurrently using parallel MCTS explorations that are 144 synced after each iteration. To achieve this we perform node expansions in the respective search 145 spaces independently but condition reward calculation and insight generation on a running optimum 146 KB state. Each iteration of EG-MCTS can be broken down two phases: 147

Expand Phase: During this stage, for each search tree, we pick the *best* state s^* and expand 148 it concurrently. To perform this expansion the KB $\mathcal{D}(s^*)$ is constructed by including the *current* 149 document $d_k(s^*)$ and the best (oracle) documents $\{d_i^*\}_{i\neq t}$ for all other positions. Thus, the KB at 150 iteration t, $0 \le t \le E$ is defined as: 151

$$\mathcal{D}_t = d_t \cup d^*_{i:i \neq t} \tag{5}$$

We use this KB $\mathcal{D}(s_i)$ to generate trajectory rollouts which are consumed by the ReflAgent to 152 generate insights. We then use the IntegAgent to generate various updated variants of d_k^* e.g., 153 $d_k(s_i), ..., d_k(s_i)$, where $0 \le i \le E$ and $0 \le j \le E$. We then estimate a reward R for each of these 154 newly generate states and update rewards of parent states using backpropagation. 155

Gather Phase: During this stage, the current best states from each document's MCTS tree are 156 gathered together and distributed to every MCTS tree for reward calculation. This is important to 157 1. Estimate rewards for each expanded state, and 2. Generate new insights for further node expansion. 158

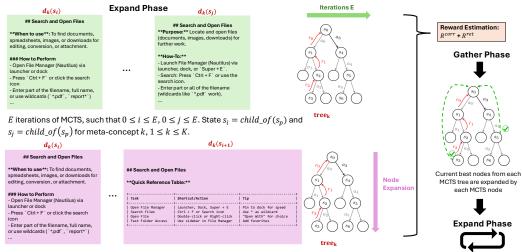
3.4 Reward-Guided Optimization

159

This section describes BREW's joint reward and loss optimization for learning an optimal KB. 160

Reward Objective: Each document state is rewarded based on two complementary criteria: (i) how 161 well the current document contributes to accurate downstream reasoning, and (ii) how retrievable 162 it is in the context of a growing KB. Formally, the total reward at time step t is defined as:

$$R_t = \lambda_{\text{corr}} \cdot R_t^{\text{corr}} + \lambda_{\text{ret}} \cdot R_t^{\text{ret}}$$
(6)



Node expansion at state s_i for meta-concept k, $1 \le k \le K$.

Figure 2: Illustration of BREW's KB optimization process using Expand-and-Gather MCTS with OSWorld examples. In the **Expand Phase**, for each document k, we sample the best node from tree $_k$ using UCT and perfrom node expansion. Node rewards are estimated based on correctness and retrievability. In the **Gather Phase**, the current best nodes from each tree are gathered per node, and the objective function is optimized. The process is repeated during the next iteration of KB refinement.

- where R_t^{corr} is the **correctness reward**, R_t^{ret} is the **retrieval reward**, and $\lambda_{\text{corr}}, \lambda_{\text{ret}} \in [0, 1]$ are scalar weights with $\lambda_{\text{corr}} + \lambda_{\text{ret}} = 1$.
- 166 **Correctness Reward:** The correctness reward R_t^{corr} evaluates the accuracy of the agent's output over a held-out query set Q, when reasoning over the current KB \mathcal{D}_t . It is defined as:

$$R^{\text{corr}}(d_t|\mathcal{D}_t) = \frac{1}{|\mathcal{Q}|} \sum_{q \in \mathcal{Q}} \text{Eval}_{\text{task}}(q, \text{agent} \oplus \mathcal{D}_t)$$
 (7)

- where $\text{Eval}_{\text{task}}$ is a task-specific evaluation function (e.g., question-answering accuracy, entailment correctness), and $\text{agent} \oplus \mathcal{D}_t$ denotes the agent acting over the hybrid KB.
- Retrieval Reward: The retrieval reward R_t^{ret} measures how effectively the current document d_t can be retrieved from the current KB \mathcal{D}_t . For a held-out query set \mathcal{Q} , it is computed using the mean reciprocal rank (MRR):

$$R^{\text{ret}}(d_t|\mathcal{D}_t) = \frac{1}{|\mathcal{Q}|} \sum_{q \in \mathcal{Q}} MRR_q(d_t, \mathcal{D}_t)$$
(8)

This encourages documents that are not only helpful in reasoning but also easily retrievable via the retrieval model over \mathcal{D}_t .

4 Experimental Setup

175

- Datasets We evaluate BREW on three diverse benchmarks testing different aspects of interactive agent capabilities: OSWORLD for computer-use automation [27], τ^2 -Bench for tool use [5], and SPREADSHEETBENCH for data manipulation [18].
- 1. **OSWorld:** This benchmark tests multimodal agents on real-world computer tasks across 10 applications. We use *GTA1-7B*, a state-of-the-art computer-use agents with BREW. Tasks are evaluated using 134 custom scripts that verify final application states.
- 182 2. τ^2 -Bench: This benchmark evaluates conversational agents on multi-turn tool-use scenarios across
 183 *Telecom*, *Retail*, and *Airline* domains. We test o4-mini-based tool-calling agent, constructing
 184 BREW KBs for every domain.

3. **SpreadsheetBench:** This benchmark evaluates agents on real-world spreadsheet manipulation, spanning both cell-level and sheet-level tasks. It contains 912 authentic user instructions paired with 2,729 test cases (3 per instruction), sourced from Excel forums and blogs. Spreadsheets include diverse formats with multi-table sheets (35.7%) and non-standard tables (42.7%). We test o4-mini using a Python tool-calling agent, and enhance it with by adding an embedding based Retrieval over the BREW KB generated over a small held-out train set of 30 samples.

Baselines We compare BREW against two widely used experiential memory approaches, *Cognee*¹ and *Agent-Mem* [30], both of which serve as established baselines for AI memory evaluation. Cognee is an open-source AI memory engine that employs a graph-plus-vector memory architecture through an Extract–Connect–Learn pipeline, enabling agents to construct cross-document and cross-context connections entirely from previously available trajectories. In contrast, Agent-Mem provides a scalable memory layer for dynamically extracting and retrieving information from conversational data, with enhanced variants incorporating graph-based memory representations. While Cognee primarily emphasizes cross-document relational reasoning, Agent-Mem focuses on scalable personalization for conversational agents.

Other Experimental Configs: For all experiments, we use GPT-4.1-2025-04-14 as the base LLM with expansion width e=3, max depth k=3, and balanced reward weights $\lambda_{\rm corr}=\lambda_{\rm ret}=0.5$. During MCTS node selection, we use the UCT [15] for balancing exploration and exploitation Full experimental details are provided in the Appendix.

5 Analysis & Discussion

In this section, we present findings from our evaluation of BREW. For more details on qualitative insights and discussion you may refer to the supplementary material.

Variations with State Search Strategy BREW performs a search across possible KB states using MCTS. We compare different state search strategies to determine the relative trade-offs:

- 1. Iterative Refinement: In this strategy we generate one version of each document to generate an initial KB, followed by a round of evals. We then use the aggregator agent to refine the documents over the newly learned insights. We repeat this step multiple times up to a maximum number of refinements. Note that in contrast to MCTS, in this strategy we do not perform node expansions and rather explore a path in the search tree.
- 2. *Greedy Search*: In this strategy we greedily pick the best state during each node expansion and only explore the sub-tree within it. This is in contrast to MCTS where, we explore different states using the UCT algorithm that balances exploration and exploitation.

Table 1 presents how MCTS achieves consistent performance gains across all benchmarks. These represent 1-5% improvements over alternative search strategies across tasks. Iterative refinement's poor performance reveals core limitations in the integrator agent feedback incorporation- which can be attributed to inherent stochasticity in LLMs. This makes state exploration especially important for textual optimization tasks like ours. We present a detailed analysis on how varying MCTS parameters result in different final states in appendix.

5.1 Trends across Sub-Tasks

BREW learns recipes from sub-trajectories in OSWorld. Figure 3 shows that BREW(BREW) improves success rates in 5 out of 10 OSWorld categories, achieving absolute gains of 4–16% while maintaining performance parity in the remaining categories (Chrome, Gimp, LibreOffice Calc, LibreOffice Impress, OS). The largest improvements appear in text-processing applications (LibreOffice Writer: $14\% \rightarrow 24\%$, Thunderbird: $38\% \rightarrow 54\%$) and multimedia tools (VLC: $20\% \rightarrow 27\%$), with moderate gains in multi-application and development environments. Even in settings with limited improvements in task correctness, BREWconsistently reduces execution length by 14–23 steps, highlighting more efficient planning.

¹github.com/topoteretes/cognee

Method	OSWorld GTA1-7B	$ au^2$ Bench o4-mini	SpreadsheetBench o4-mini
Baseline	44.20	56.63	44.30
Cognee	46.70	57.71	42.10
Agent-Mem	43.83	52.69	42.00
BREW-Iterative	46.13	57.34	42.98
BREW-Greedy	45.55	59.14	45.94
BREW-MCTS	47.56	59.14	46.80

Table 1: Comparison of models under different evaluation setups, including Baseline model and BREW augmented model. We report task success rate for OSWorld, ratio of independent tasks that succeeded for τ^2 Bench, and the 1st test case pass rate for SpreadsheetBench.

OSWorld: Success Rate Comparison and Efficiency Gains

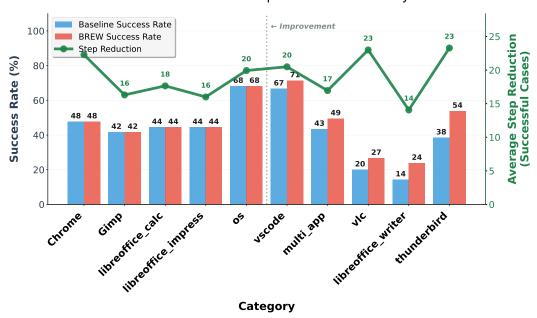


Figure 3: The bar plot represents the category-wise success rate over various tasks in the OSWorld dataset over the GTA1-agent, whereas the line plot demonstrates the reduction in the number of steps for the successful cases. Note that even in scenarios where the KB doesn't help increase the success rate, it significantly reduces the number of steps needed to succeed.

This pattern suggests that BREW's architectural enhancements are particularly effective for tasks requiring complex sequential reasoning and inter-application coordination, while preserving baseline robustness in domains constrained by intrinsic task complexity.

A qualitative analysis of the knowledge bases (KBs) constructed by BREWfurther supports this finding. We observe that BREWcaptures and represents sub-trajectory characteristics in *natural language*, including application shortcuts, standard operating procedures, and strategies for localizing UI elements. Since many UI tasks share common sub-trajectories, this representation facilitates knowledge transfer across tasks within the same application. Moreover, BREWsubstantially reduces reliance on granular UI interactions: while the baseline GTA1 model executes approximately 19,000 clicks and 17,821 keyboard actions, BREWsignificantly decreases this interaction complexity.

BREW learns aggressive resolution strategies for $\tau^2 - Bench$ To evaluate robustness of BREW, we analyzed the distribution of failure modes across the τ^2 -retail dataset, focusing on four key error categories: Wrong Argument, Wrong Info, Wrong Decision, and Partially Resolve. Figure 4 presents a comparative chart for the baseline, BREW, Cognee and Agent-Mem[30].

Overall, BREW demonstrated consistent improvements across most error types compared to the baseline and competing approaches. Specifically, BREW showed a **notable reduction in "Wrong Argument" and "Wrong Decision" errors**, indicating that it was better at capturing logical dependencies in retail dialogues and making accurate decisions.

Interestingly, *Partially Resolve* errors were slightly higher for BREW than for Cognee, likely because BREW attempted more aggressive resolution strategies that occasionally failed to fully satisfy user queries. Cognee appears to capture *richer factual details* given its relatively lower *Wrong Info* errors, whereas Agent-Mem excels in *tracking conversation state and decision accuracy*, as reflected in its reduced *Wrong Decision* failures.

Improvements in Task Efficiency We observe that overall, BREWenables agents to come to a correct response quicker.

OSworld. Figure 3 demonstrates that BREW enables GTA1 to complete tasks more efficiently. Compared to the baseline GTA1 model's average of \sim 75 steps, the BREW-augmented model completes tasks 14% faster with an average of \sim 64 steps. Analyzing performance by outcome reveals that while step counts remain unchanged for failed cases, successful completions show a substantial 39% (rel.) reduction in execution steps, indicating improved planning efficiency for achievable tasks.

 τ^2 Bench. Similarly, BREW reduces average conversation turns from 29.47 to 28.43 (-3.5%), while maintaining consistent step reductions across categories. Step reductions average 1.7 steps for Retail and Telecom, but 3.1 steps for Airline, indicating greater efficiency gains in complex domains. Qualitative analysis seconds these numbers showing how knowledge base integration enables more direct task completion paths and improved planning quality, though multi-turn interactions remain necessary for complex sub-tasks.

SpreadsheetBench. While we observe a slight increase in the number of turns across the entire benchmark suite $(4.5 \rightarrow 5.4)$ in the case of the baseline versus BREW, an interesting pattern

emerges in more than 82% of the cases the baseline and the BREW appended agent performs

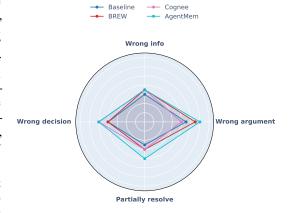


Figure 4: Distribution of errors in τ^2 Bench Retail

similarly with similar turn consumption. BREW leads to an improvement in 12% of the cases where the KB is able to address gaps in the baseline technique to enable the agent to go exploring further leading to positive outcomes with an average of 1 step increase in the interactions.

6 Conclusions

In this work, we explored an alternative approach to agent optimization by focusing on experiential knowledge retention rather than direct model fine-tuning. We introduced BREW, a framework that aims to construct and refine a structured, interpretable knowledge base from past agent interactions. By decomposing agent memory into concept-level documents and applying a state-search optimization strategy, BREW provides a modular and transparent substrate for memory formation. Our evaluations across OSWorld and τ^2 Bench benchmarks suggest that such structured memory can support measurable improvements in task success and efficiency, while maintaining manageable computational costs. Although the observed gains are promising, we recognize that BREW's effectiveness is influenced by the quality and coverage of its training data. Future work could explore more adaptive and domain-general memory refinement techniques, as well as tighter integrations with ongoing agent planning. Ultimately, we hope this study encourages further investigation into more interpretable, memory-driven approaches to language agent development—especially in real-world environments where long-term consistency and adaptability are essential.

7 References

- Lakshya A. Agrawal, Shangyin Tan, Dilara Soylu, Noah Ziems, Rishi Khare, Krista Opsahl-Ong, Arnav Singhvi, Herumb Shandilya, Michael J. Ryan, Meng Jiang, Christopher Potts, Koushik
 Sen, Alexandros G. Dimakis, Ion Stoica, Dan Klein, Matei Zaharia, and Omar Khattab.
 Gepa: Reflective prompt evolution can outperform reinforcement learning. arXiv preprint arXiv:2507.19457, July 2025.
- 203 [2] Anthropic. Introducing computer use, a new Claude 3.5 Sonnet, and Claude 3.5 Haiku, October 2024. URL https://www.anthropic.com/news/3-5-models-and-computer-use. Accessed: 2025.
- Yasharth Bajpai, Bhavya Chopra, Param Biyani, Cagri Aslan, Dustin Coleman, Sumit Gulwani,
 Chris Parnin, Arjun Radhakrishna, and Gustavo Soares. Let's fix this together: Conversational
 debugging with github copilot. In 2024 IEEE Symposium on Visual Languages and Human Centric Computing (VL/HCC), pages 1–12, 2024. doi: 10.1109/VL/HCC60511.2024.00011.
- [4] Victor Barres, Honghua Dong, Soham Ray, Xujie Si, and Karthik Narasimhan. τ^2 -bench: Evaluating conversational agents in a dual-control environment, 2025. URL https://arxiv.org/abs/2506.07982.
- 5313 [5] Victor Barres, Honghua Dong, Soham Ray, Xujie Si, and Karthik Narasimhan. τ^2 -bench: Evaluating conversational agents in a dual-control environment, 2025. URL https://arxiv.org/abs/2506.07982.
- [6] Param Biyani, Yasharth Bajpai, Arjun Radhakrishna, Gustavo Soares, and Sumit Gulwani.
 Rubicon: Rubric-based evaluation of domain-specific human ai conversations. In *Proceedings*of the 1st ACM International Conference on AI-Powered Software, AIware 2024, page 161–169,
 New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400706851. doi:
 10.1145/3664646.3664778. URL https://doi.org/10.1145/3664646.3664778.
- [7] 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.
- [8] Rémi Coulom. Efficient selectivity and backup operators in monte-carlo tree search. In *Proceedings of the 5th International Conference on Computers and Games (CG 2006)*, pages 72–83. Springer, 2006. doi: 10.1007/978-3-540-75538-8_7.
- [9] Lutfi Eren Erdogan, Nicholas Lee, Sehoon Kim, Suhong Moon, Hiroki Furuta, Gopala Anumanchipalli, Kurt Keutzer, and Amir Gholami. Plan-and-act: Improving planning of agents for long-horizon tasks. *The Forty-Second International Conference on Machine Learning*, 2025.
- [10] Priyanshu Gupta, Shashank Kirtania, Ananya Singha, Sumit Gulwani, Arjun Radhakrishna, Gustavo Soares, and Sherry Shi. MetaReflection: Learning instructions for language agents using past reflections. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 8369–8385, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.477. URL https://aclanthology.org/2024.emnlp-main.477/.
- [11] Maryam Hashemzadeh, Elias Stengel-Eskin, Sarath Chandar, and Marc-Alexandre Cote. Sub goal distillation: A method to improve small language agents, 2024. URL https://arxiv.
 org/abs/2405.02749.
- 340 [12] Yuanzhe Hu, Yu Wang, and Julian McAuley. Evaluating memory in llm agents via incremental multi-turn interactions, 2025. URL https://arxiv.org/abs/2507.05257.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=VTF8yNQM66.

- [14] Zaid Khan, Elias Stengel-Eskin, Archiki Prasad, Jaemin Cho, and Mohit Bansal. Executable 346 functional abstractions: Inferring generative programs for advanced math problems. 2025. 347
- [15] Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In Johannes 348 Fürnkranz, Tobias Scheffer, and Myra Spiliopoulou, editors, Machine Learning: ECML 2006, 349 pages 282-293, Berlin, Heidelberg, 2006. Springer Berlin Heidelberg. ISBN 978-3-540-46056-350 351
- [16] Xinzhe Li. A review of prominent paradigms for llm-based agents: Tool use, planning (including 352 rag), and feedback learning. In Proceedings of the 31st International Conference on Compu-353 tational Linguistics (COLING), pages 9760-9779, Abu Dhabi, UAE, 2025. Association for 354 Computational Linguistics. URL https://aclanthology.org/2025.coling-main.652. 355
- [17] Michael L. Littman. An optimization-based categorization of reinforcement learning environ-356 ments. 1993. URL https://api.semanticscholar.org/CorpusID:17988064. 357
- [18] Zeyao Ma, Bohan Zhang, Jing Zhang, Jifan Yu, Xiaokang Zhang, Xiaohan Zhang, Sijia Luo, 358 Xi Wang, and Jie Tang. Spreadsheetbench: Towards challenging real world spreadsheet 359 manipulation. Advances in Neural Information Processing Systems, 37:94871–94908, 2024. 360
- Introducing Operator, January 2025. URL https://openai.com/index/ 361 introducing-operator/. Accessed: 2025. 362
- [20] Charles Packer, Sarah Wooders, Kevin Lin, Vivian Fang, Shishir G. Patil, Ion Stoica, and 363 Joseph E. Gonzalez. Memgpt: Towards Ilms as operating systems, 2024. URL https: 364 //arxiv.org/abs/2310.08560. 365
- [21] Yujia Qin, Yining Ye, Junjie Fang, Haoming Wang, Shihao Liang, Shizuo Tian, Junda Zhang, 366 Jiahao Li, Yunxin Li, Shijue Huang, et al. Ui-tars: Pioneering automated gui interaction with 367 native agents. arXiv preprint arXiv:2501.12326, 2025. 368
- [22] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and 369 Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model, 370 2024. URL https://arxiv.org/abs/2305.18290. 371
- [23] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal 372 policy optimization algorithms. In Proceedings of the 34th International Conference on Machine 373 Learning (ICML 2017), 2017. URL https://arxiv.org/abs/1707.06347. 374
- 375 [24] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of 376 mathematical reasoning in open language models, 2024. URL https://arxiv.org/abs/ 377 2402.03300. 378
- [25] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 379 Reflexion: Language agents with verbal reinforcement learning. In Proceedings of the 380 37th Conference on Neural Information Processing Systems (NeurIPS 2023), New Orleans, 381 LA, USA, 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/ 382 hash/1b44b878bb782e6954cd888628510e90-Abstract-Conference.html. 383
- [26] Yu Wang, Chi Han, Tongtong Wu, Xiaoxin He, Wangchunshu Zhou, Nafis Sadeq, Xiusi Chen, 384 Zexue He, Wei Wang, Gholamreza Haffari, Heng Ji, and Julian McAuley. Towards lifespan 385 cognitive systems, 2025. URL https://arxiv.org/abs/2409.13265. 386
- [27] Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, 387 Toh Jing Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu, Yiheng Xu, Shuyan 388 Zhou, Silvio Savarese, Caiming Xiong, Victor Zhong, and Tao Yu. Osworld: Benchmarking 389 multimodal agents for open-ended tasks in real computer environments. In A. Globerson, 390 L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors, Advances in 391 Neural Information Processing Systems, volume 37, pages 52040–52094. Curran Associates, 392 Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/ 393 $5d413e48f84dc61244b6be550f1cd8f5-Paper-Datasets_and_Benchmarks_Track.$ 394

pdf. 395

- Ran Xu, Yuchen Zhuang, Yue Yu, Haoyu Wang, Wenqi Shi, and Carl Yang. Rag in the wild: On the (in)effectiveness of llms with mixture-of-knowledge retrieval augmentation. *arXiv preprint* arXiv:2507.20059, 2025.
- 399 [29] 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.
- 401 [30] Wujiang Xu, Kai Mei, Hang Gao, Juntao Tan, Zujie Liang, and Yongfeng Zhang. A-mem:
 402 Agentic memory for llm agents, 2025. URL https://arxiv.org/abs/2502.12110.
- 403 [31] Hui Yang, Sifu Yue, and Yunzhong He. Auto-gpt for online decision making: Benchmarks and
 404 additional opinions. *arXiv preprint arXiv:2306.02224*, 2023. doi: 10.48550/arXiv.2306.02224.
 405 URL https://doi.org/10.48550/arXiv.2306.02224.
- John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik
 Narasimhan, and Ofir Press. Swe-agent: Agent-computer interfaces enable automated software
 engineering. Advances in Neural Information Processing Systems, 37:50528–50652, 2024.
- 409 [33] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan
 410 Cao. React: Synergizing reasoning and acting in language models. In *Proceedings of the*411 11th International Conference on Learning Representations (ICLR 2023), 2023. URL https:
 412 //openreview.net/forum?id=WE_vluYUL-X.
- Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. τ -bench: A benchmark for tool-agent-user interaction in real-world domains. In *NeurIPS (Workshops)*, 2024. State-of-the-art agents (e.g. GPT-40) succeed on <50
- Yuyan Zhou, Liang Song, Bingning Wang, and Weipeng Chen. Metagpt: Merging large language models using model exclusive task arithmetic. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1711–1724, Miami, Florida, USA, 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024. emnlp-main.102.

A Appendix

2 A.1 Details of the BREWAlgorithm

We provide pseudocode for the core components of BREW, aligning with the stages introduced in Section 3. Each algorithm plays a distinct role in constructing, organizing, or refining the knowledge base over iterative interactions. GenerateInsights (Alg. 2) produces concept-aligned insights from annotated rollouts using ReflAgent. DeduplicateConcepts (Alg. 3) clusters semantically overlapping concepts into a compact meta-concept set. Integagent incrementally builds and updates per-concept documents using newly generated insights. Finally, Expanding (Alg. 4) performs MCTS-guided expansions to explore improved document variants, while Evaluate (Alg. 5) scores candidate KB states using correctness and retrieval-based rewards.

We specify the IntegAgent prompt below:

BREW Integrator Prompt

432 433

```
434
        # Enhanced Documentation Editor Prompt
435
        You are a meticulous documentation-level editor specializing in comprehensive
436
            technical reference materials. You will be given a list of topic nodes,
437
             each containing structured information that must be preserved and enhanced
438
            with maximum detail retention.
439
440
        ## Input Structure Analysis
441
        Each node contains:
442
443
        - **Title**: The primary topic identifier
444
        - **Context**: Background information and conceptual foundation
        - **How to Use**: Step-by-step instructions, commands, flags, parameters, and
445
446
            implementation details
        - **When to Use**: Specific scenarios, conditions, and decision criteria
447
        - **Best Practices**: Expert recommendations, optimization techniques, and
448
449
             common pitfalls to avoid
450
451
        ## Detailed Processing Requirements
452
        ### 1. Information Preservation (Zero Loss Policy)
453
        - **Preserve every technical detail**: All command-line flags, parameter values,
454
             configuration options, file paths, URLs, version numbers, and exact syntax
455
        - **Maintain all examples**: Keep every code snippet, sample input/output, file
456
457
            names, directory structures, and command sequences exactly as provided
458
        - **Retain contextual nuances**: Preserve qualifying language like "typically,"
             "usually," "in most cases," "when available," and conditional statements
459
          **Keep quantitative data**: Preserve all numbers, measurements, timeframes,
460
            limits, thresholds, and statistical information
461
462
        - **Maintain cross-references**: Keep all mentions of related tools,
             dependencies, prerequisites, and interconnected concepts
463
464
        ### 2. Enhanced Detail Extraction
465
          **Expand abbreviations**: When encountering shortened forms, expand them
466
            naturally while preserving the original
467
        - **Surface implicit knowledge**: Make obvious assumptions explicit (e.g., "this
468
469
             requires root permissions," "assumes default configuration")
        - **Clarify relationships**: Explicitly describe how different components,
470
471
             options, or steps relate to each other
        - **Highlight edge cases**: Emphasize special conditions, exceptions, or unusual
472
             scenarios mentioned in the source
473
        - **Elaborate on consequences**: When the source mentions outcomes, expand on
474
475
            both success and failure scenarios
476
        ### 3. Prose Transformation Guidelines
477
478
          **Bullet integration**: Transform each bullet point into 1-3 complete
            sentences that naturally flow together
479
```

```
- **Technical precision**: Use precise technical vocabulary while maintaining
480
             readability
481
         - **Logical flow**: Organize information within each section to follow a logical
482
              sequence (setup \rightarrowexecution \rightarrowverification)
483
         - **Contextual embedding**: Weave code snippets and technical terms seamlessly
484
             into narrative sentences
485
         - **Comprehensive coverage**: Ensure every sub-bullet, nested item, and
486
             parenthetical note becomes part of the prose
487
488
489
        ### 4. Structural Requirements
490
         - **Heading hierarchy**: Use '# Title' for each node's main heading
        - **Section order**: Maintain Context \rightarrowHow to Use \rightarrowWhen to Use \rightarrowBest
491
             Practices sequence
492
         - **Paragraph organization**: Create substantial paragraphs (3-6 sentences)
493
             rather than brief statements
494
495
         - **Transition quality**: Craft smooth bridges between sections and between
             different nodes
496
         - **Code formatting**: Preserve all inline code with backticks and maintain
497
             proper formatting for code blocks
498
499
        ### 5. Quality Assurance Checklist
500
        Before finalizing, verify:
501
        - [ ] Every piece of source information appears in the output
502
        - [ ] All technical specifications, parameters, and examples are intact
503
        - [ ] Code snippets maintain their exact syntax and formatting
504
         - [ ] Prose flows naturally without choppy or fragmented sentences
505
         - [ ] Each section provides comprehensive coverage of its topic area
506
        - [ ] Cross-references and dependencies are clearly explained
507
        - [ ] No section labels or formatting artifacts remain in the prose
508
509
        ## Output Specifications
510
        Generate a single, cohesive markdown document that reads as authoritative
511
             technical documentation. The result should be comprehensive enough that a
512
             reader could successfully implement the described tools or techniques using
513
              only the information provided, without referring back to the original
514
515
             nodes.
516
517
518
        **Input Nodes:**
519
        <NODES>
520
        {node_list}
521
        </NODES>
522
523
524
525
        Now, produce the aggregated markdown reference sheet with maximum detail
526
             preservation and enhanced clarity.
528
```

Algorithm 2 GenerateInsights: Extract behavioral insights from trajectories

```
Require: Queries \mathcal{Q}, KB \mathcal{D}, rubrics
Ensure: Concept-insight pairs \mathcal{B}
 1: \mathcal{B} \leftarrow \emptyset
 2: for each query q \in \mathcal{Q} do
 3:
           \tau \leftarrow \text{LLM}(q, \mathcal{D})
                                                                                                               4:
           label \leftarrow Grade(\tau)
                                                                                                                     (c, i) \leftarrow \text{REFLAGENT}(\tau, \text{rubrics}, \text{label})
 5:
           \mathcal{B} \leftarrow \mathcal{B} \cup \{(c, i, q)\}

    Store with source query

 7: end for
 8: return \mathcal{B}
```

```
Algorithm 3 DeduplicateConcepts: Cluster similar concepts and map queries
```

```
Require: Concept-insight-query triples \mathcal{B}
Ensure: Meta-concepts \mathcal{K} with mapped queries and insights

1: Extract all concepts from \mathcal{B}

2: Embed and cluster concepts by similarity

3: \mathcal{K} \leftarrow cluster representatives

4: for each k \in \mathcal{K} do

5: \mathcal{Q}_k^{\text{train}} \leftarrow {training queries that contributed insights to k}

6: \mathcal{Q}_k^{\text{eval}} \leftarrow {held-out queries relevant to k}

7: \mathcal{I}_k \leftarrow {all insights mapped to concept k}

8: end for

9: return \mathcal{K} with associated queries and insights
```

Algorithm 4 ExpandNode: Generate and evaluate new document variants

```
Require: Node s, concept k, candidates h, current KB \mathcal{D}_{current}, best docs \mathcal{D}_{best}, tree
Ensure: Updated tree with new evaluated nodes
                                                                  ▷ Generate new insights from concept-relevant queries
 1:
 2: \mathcal{B}_{new} \leftarrow \varnothing
 3: for query q \in \mathcal{Q}_k^{\text{train}} do 4: \tau \leftarrow \text{LLM}(q, \mathcal{D}_{\text{current}})
 5:
           (c, i) \leftarrow \text{ANNOTATE}(\tau, \text{rubrics}, \cdot)
           if c maps to k then
 6:
 7:
                 \mathcal{B}_{\text{new}} \leftarrow \mathcal{B}_{\text{new}} \cup \{i\}
 8:
           end if
 9: end for
10:
                                                                                11: for j = 1 to h do
           d_{k,j} \leftarrow \text{INTEGAGENT}(k, \mathcal{I}_k \cup \mathcal{B}_{\text{new}}, d_k^s)
12:
                                                     ▶ Evaluate using hybrid KB with best docs from other concepts
13:
           \mathcal{D}_{\text{hybrid}} \leftarrow \{d_{k,j}\} \cup \{d_{k'} \in \mathcal{D}_{\text{best}} : k' \neq k\}
14:
           R_{k,j} \leftarrow \text{EVALUATE}(d_{k,j}, \mathcal{D}_{\text{hybrid}}, \mathcal{Q}_k^{\text{eval}})
15:
16:

    ▷ Add to tree and backpropagate

           Add (d_{k,j}, R_{k,j}) as child of s in tree
17:
           Backpropagate R_{k,i} from new node to root
18:
19: end for
```

Algorithm 5 Evaluate: Score document using held-out queries

```
Require: Document d_k, hybrid KB \mathcal{D}_{\text{hybrid}}, eval queries \mathcal{Q}_k^{\text{eval}}
Ensure: Reward score R

1: R^{\text{corr}} \leftarrow 0
2: R^{\text{ret}} \leftarrow 0
3: for each q \in \mathcal{Q}_k^{\text{eval}} do
4: R^{\text{corr}} \leftarrow R^{\text{corr}} + \text{EVAL}(q, \text{agent} \oplus \mathcal{D}_{\text{hybrid}})
5: R^{\text{ret}} \leftarrow R^{\text{ret}} + \text{MRR}(d_k, q, \mathcal{D}_{\text{hybrid}})
6: end for
7: R^{\text{corr}} \leftarrow \frac{R^{\text{corr}}}{|\mathcal{Q}_k^{\text{eval}}|}
8: R^{\text{ret}} \leftarrow \frac{R^{\text{ret}}}{|\mathcal{Q}_k^{\text{eval}}|}
9: return \lambda_{\text{corr}} \cdot R^{\text{corr}} + \lambda_{\text{ret}} \cdot R^{\text{ret}}
```

529 A.2 BREW Configurations

Base LLM Configuration For all BREWalgorithm steps, we use the OpenAI GPT-4.1-2025-04-14 model as the underlying language model. To balance exploration and stability, we set the temperature

- to 0.7 for the IntegAgent component to encourage diversity in sampled completions, while all other
- calls use a temperature of 0.1 for deterministic behavior. The search process employs an expansion
- width of e=3, a maximum search depth of k=3, and a maximum of n=10 iterations. Reward
- signals are weighted equally across correctness and retrieval relevance, with $\lambda_{corr} = \lambda_{ret} = 0.5$.

536 A.3 Baseline Methods

- 537 We compare BREWagainst two common reasoning baselines. Step-Back Prompting encourages
- backward reasoning by guiding the model to work from the final task objective back to the initial
- actions. In-Context Learning augments the input prompt with successful trajectories from related
- tasks, enabling the model to benefit from relevant prior examples without additional fine-tuning.

541 A.4 Benchmark Specifications

542 A.4.1 OSWorld: Computer-Use Automation

- Dataset Overview OSWorld [27] comprises 369 real-world computer-use tasks spanning 10 distinct
- applications. The benchmark is divided into train and test sets, with the distribution of tasks across
- domains shown in Table 2.
- Agent Specifications The UI-Tars-7B variant is a 7B-parameter multimodal transformer fine-tuned
- for graphical user interface understanding. It operates over an action space of PyAutoGUI commands
- (e.g., click, type, and key presses). The agent integrates a retrieval module that queries a task-relevant
- knowledge base using the user-provided description, with the top three retrieved items added to the
- system prompt. Inputs to the model consist of a screenshot of the active UI paired with the natural
- 551 language task description.
- 552 The GTA1-7B configuration adopts a two-agent architecture, consisting of a planner and a grounding
- module. The planner (GTA-1-7B) generates the high-level action sequence, while the grounding
- module (OpenAI O3) verifies and refines each action before execution. Knowledge retrieval is
- incorporated differently for each component: the planner performs a single retrieval at the start
- of execution, which is persisted in its prompt, whereas the grounding module performs dynamic
- retrievals at each verification step.
- Evaluation Protocol Evaluation uses 134 task-specific scripts designed for automated verification.
- 559 Success criteria include file state checks (e.g., validating .xlsx or .docx outputs), UI element
- validation to confirm correct interaction, and process completion checks to ensure that the intended
- automation sequence was executed successfully.

562 **A.4.2** τ^2 -Bench: Interactive Tool Usage

- Dataset Overview τ^2 -Bench [5] extends τ -Bench by introducing bidirectional tool-calling capa-
- bilities. The dataset covers multiple service-oriented domains, with domain-level task distributions
- summarized in Table 3.
- 566 **Domain Characteristics** The benchmark spans several domains with distinct task characteristics.
- The Telecom domain focuses on connectivity troubleshooting, plan modifications, and service
- activation workflows. The Retail domain includes order processing, return handling, and inventory
- queries. The Airline domain emphasizes booking modifications and policy-compliant rescheduling
- 570 scenarios.
- 571 Interaction Settings Two interaction modes are defined. In Easy mode, a human proxy (imple-
- mented via GPT-4.1) provides detailed guidance to the agent. The knowledge base is built exclusively
- 573 from Easy mode trajectories, ensuring high-quality demonstrations for learning. In Hard mode,
- human intervention is minimized. The knowledge base combines both Easy and Hard trajectories,
- testing the agent's robustness to underspecified or noisy instructions.
- 576 Evaluation Criteria Task success is measured using domain-specific verification procedures. These
- include database state checks to validate final outcomes, status checks for confirming service or
- connection state, natural language verification to ensure correct confirmation statements appear in
- 579 dialogue, and action matching to confirm that all required steps are completed. Each domain uses a
- tailored subset of these checks (e.g., Telecom relies primarily on status checks).

Domain	Test	Train
Calc	45	2
Chrome	44	2
Writer	21	2
Gimp	24	2
Impress	45	2
Os	22	2
Thunderbird	13	2
Multi-apps	99	2
VLC	15	2
VSCode	21	2
Total	349	20

Table 2: Test and Train samples across different domains in OSWorld.

Domain	Test	Train
Telecom	105	7
Retail	105	7
Airline	44	6
Total	254	20

Table 3: Task-wise breakdown for τ^2 -Bench with assumed 2-shot training samples per domain.

581 Domain Characteristics

583

584

586

587

588

589

591

592

593

594

595

596

597

598

599

600

601

602

- **Telecom:** Connectivity issues, plan management, service activation
 - **Retail:** Order processing, returns, inventory queries
 - Airline: Booking modifications, policy-compliant rescheduling

585 **Evaluation Criteria** Task success determined by:

- Database Checks: Final state verification
 - Status Checks: Service/connection state validation
 - NL Checks: Confirmation statements in dialogue
 - Action Matching: Required action sequence completion

Note: Each domain uses specific check combinations (e.g., Telecom uses only status checks).

A.4.3 SpreadsheetBench: Real-World Spreadsheet Manipulation

Dataset Overview SpreadsheetBench [18] consists of 912 instructions collected from four major Excel forums and blogs. Each instruction is paired with spreadsheets reflecting authentic, complex user scenarios, often containing multiple tables and non-standard relational structures. The dataset totals 2,729 test cases, averaging three per instruction. A breakdown of cell-level and sheet-level manipulations is shown in Table 4.

Task Settings The benchmark defines two dimensions of evaluation:

- **Granularity:** Instructions involve either *cell-level* manipulations (specific ranges such as D2:D6) or *sheet-level* manipulations (entire tables or multi-sheet updates).
- **Evaluation:** Performance is measured using an Online Judge (OJ)-style protocol. The *soft* setting (IOI-style) awards partial credit when only some test cases are solved, while the *hard* setting (ICPC-style) requires solutions to succeed on all test cases.

603 **Agent Configuration** We evaluate

texttto4-mini using a function-calling agent connected to a single Python execution tool. The

agent translates natural language instructions into Python code for spreadsheet manipulation (e.g., modifying cells, applying formulas, restructuring tables). After each tool call, all formulas in the spreadsheet are recalculated to ensure consistency before proceeding to the next step. This setup provides a controlled environment to assess reasoning, code generation, and execution robustness across diverse spreadsheet tasks.

Granularity	Instructions	Test Cases
Cell-Level	329	986
Sheet-Level	583	1,743
Total	912	2,729

Table 4: Cell-level vs. sheet-level distribution in SpreadsheetBench.

610 A.5 KB Construction and Retrieval Details

611 Training Data Collection

- OSWorld: 20 successful trajectories (2 per application domain) and 10 for evals.
- τ^2 -Bench: 20 trajectories balanced across domains and difficulty settings and 10 for evals.
- SpreadsheetBench: Uniformly sample 30 trajectories for training and 10 for evaluation.
- All numbers are reported on the remaining train set.

616 Retrieval Strategy

612

613

614

617 618

619

620

621

622

625

626

627

628

629

630

631

- **Query Formation:** For each task we take in the seed Natural Language query as the retrieval query.
 - Retrieval Count: We take top-3 documents for all the retrieval steps
- Integration Point: For SPREADSHEET ENCH and OSWorld we insert retrievals in the system prompt augmentation. For τ^2 -bench we add perfrom retrieval after each user interaction.

B Qualitative Analysis

- Exploration on MCTS parameters WE evaluate OSworld on two different MCTS parameters.
 - Increased Depth: To increase the depth we keep maximum width of the tree as 3 and depth as 10 with max number of iterations as 25. We observe that the Knowledge base over optimizes on the train set leading to a poorer performance on test set.
 - Increased Width: For increased width we reverse the parameters where depth is 3 and maximum width is 10 with max iterations 25. We observe many different styles of KBs are generated storing very similar information, these different styles lead to a varied performance on both eval and test set notifying the importance of state search.
- We report the numbers on table ??

	Baseline	max_width=3, max_depth=3	max_width=3, max_depth=10	max_width=10, max_depth=3
OSworld	44.20	47.56	43.83	49.32

Table 5: OSworld difference in MCTS parameters

C Exemplar Knowledge Bases

634

635

636

637

C.1 Knowledge base learned for OSWorld

We showcase a small part of knowledge base learned thought BREW. This demonstrate 3 major parts on which each document is aggregated. These parts discuss when to use a piece of information, why to use the information, how to use the information/tool.

```
638
         ## Search and Open Files
639
640
         **When to use**: Locating documents, spreadsheets, images, or downloads for
641
642
              editing, conversion, or attachment.
643
         ### How to Perform
644
         - Open **File Manager (Nautilus)** from launcher or system dock
645
         - Press 'Ctrl + F' or click the search icon
646
         - Enter part of filename, full name, or wildcard ('*.pdf', 'report*')
647
         - Use right-click \rightarrow **Open With** to choose the desired application
648
649
         - Use the sidebar to navigate to **Downloads**, **Documents**, or custom folders
650
651
         ### Additional Actions
         - Right-click →**Properties** to check modification date or file type
652
         - Sort results by Date, Type, or Name from the top-right dropdown
653
         - Use 'F2' to rename files inline
654
655
         ### Example
656
         - Task: "Edit the file titled 'sales_report_march.ods'"
657
658
           - Search for 'sales' in File Manager
           - Confirm '.ods' type and open with LibreOffice Calc
659
660
661
662
         ## Insert Images
663
664
         **When to use**: Adding visual elements to documents, presentations, emails, or
665
666
              templates.
667
         ### How to Perform
668
         - Navigate to **Insert →Image →From File** (in Writer, Impress, Thunderbird)
669
         - Select an image file ('.png', '.jpg', '.svg') from the file dialog
670
         - Use drag handles to resize; right-click \to**Wrap** or **Alignment** for layout
671
672
         ### Additional Actions
673
         - In GIMP: **File \rightarrowOpen as Layers** to insert image as a new layer
674
675
         - Use drag-and-drop from file manager into open document windows
         - Use **Format →Image** to apply borders, shadows, or color corrections (in
676
677
             Writer/Impress)
678
679
         - Task: "Insert the logo.png image into the title slide"
680
           - Open '.odp' file in Impress \rightarrow \texttt{Go} to Slide 1 \rightarrow \texttt{Insert} \rightarrow \texttt{Image} \rightarrow \texttt{Select} 'logo.
681
               png'
682
683
684
         . . .
685
         ## Export as PDF
686
687
         **When to use**: Required submission format
688
689
690
         ### How to Perform
         - Go to **File \rightarrowExport As PDF**
691
         - Choose output folder (usually **Documents** or **Downloads**)
692
693
         - Click **Save**, then confirm the exported file opens correctly
694
```

```
### Additional Actions
- In GIMP or Impress: choose **File →Export As**, then select '.pdf' from
format list
- Use **Save As** to preserve both editable and exported versions separately
### Example
- Task: "Export the flyer.xcf as a PDF"
- Open in GIMP →File →Export As →Rename to 'flyer.pdf' →Click Export
```

C.2 BREW Knowledge Base for τ^2 -Bench

706

BREW enable use to learn relevant information for tau bench for across the domains in a single knowledge base. This knowledge base is helpful to use relevant actions from the action pool.

```
707
        ### Additional Actions
709
710
        * Inform the user:
711
712
           - Refunds via gift card = immediate.
713
          - Refunds via other methods = -57 business days.
714
        ### Example
715
716
        * Task: "Cancel a T-shirt order placed yesterday"
717
718
          * Validate: Status is 'pending'
          * Reason: "no longer needed"
719
           * Confirm
720
721
          * Execute tool call
722
723
        # Exchange Delivered Order
724
725
726
        **When to use**:
        User wants to swap delivered items for a different variant (e.g., size or color)
727
728
729
730
         **Why to use it**:
        To fix sizing or option errors without needing a new purchase.
731
732
        ### How to Perform
733
         - Authenticate user
734
        - Confirm order status is 'delivered'
735
736
         - Get full list of exchange items
        > "Please ensure all items for exchange are listed. This step 'cant be repeated.
737
738
         - Ask for refund/payment method
739
740
        - Confirm:
          > "'Youre exchanging item X for same product, different option. Proceed?"
741
         - On confirmation:
742
           ""python
743
          request_exchange(order_id="45678", item_exchanges=[...], payment_method="
744
               paypal")
745
746
747
        ### Additional Actions
748
749
        * Mention: An email will be sent with return instructions
750
        * Validate that the new variant is from the same product
751
752
        ### Example
753
754
         * Task: "Exchange red shirt for blue in Order #45678"
755
756
           * Confirm all exchange items
757
          * Confirm payment method for difference
```

```
* Execute tool call
758
759
         ### Example
760
761
         * Task: "Show me my last 2 orders"
762
          * Authenticate
763
764
          * Retrieve and present info
765
         # Deny Unsupported Request
766
767
768
         **When to use**:
         User asks for an unsupported action (e.g., cancel processed order, exchange to
769
             different product type, help another user).
770
771
         **Why to use it**:
772
773
         To stay compliant with platform policy.
774
         ### How to Perform
775
         - Politely reject:
776
          > "'Im sorry, but I 'cant process that request. 'Its outside the allowed scope.
777
778
779
         ### Example
780
781
         * Task: "Cancel a processed order"
782
          * Respond with denial message
783
         # Transfer to Human Agent
784
785
786
         **When to use**:
         User needs help outside the 'assistants permitted capabilities.
787
788
         **Why to use it**:
789
         To ensure user gets the right help from trained staff.
790
791
         ### How to Perform
792
         - Make tool call:
793
          ""python
794
795
          transfer_to_human_agents()
796
         - Then inform user:
797
          > "YOU ARE BEING TRANSFERRED TO A HUMAN AGENT. PLEASE HOLD ON."
798
799
         ### Example
800
801
         * Task: "Delete a task"
802
803
           * Deny deletion
           * Transfer to human
804
```

C.3 BREW Knowledge Base for SpreadsheetBench

806

```
807
            Header Extraction
808
        1. Detecting Header Rows
809
810
811
        To accurately identify header rows, scan the initial region of your dataset.
             This process is crucial for mapping column information for further
812
813
             processing.
814
        Approaches:
815
         - Heuristic Checks:
816
817
        - Look for rows where all cells are strings (e.g., "Name", "Date", "Region", "
             Amount").
818
         - Identify rows with distinctive formatting such as bold text or background
819
820
             color.
821
         - Example:
```

```
| Name | Date | Region | Amount | |-----| |
822
             John | 2024-01-01| North | 100 |
823
824
        - Pattern Recognition:
        - Use regex to match typical header patterns, such as column names starting with
825
              uppercase letters.
826
        - Score candidate rows based on the likelihood of being headers.
827
        - Multi-Table Sheets:
828
        - Detect gaps, empty rows, or separators indicating a new table.
829
        - Assign a Table ID to each detected table for later reference.
830
831
832
        Edge Cases:
        - Merge multi-row headers (e.g., "Sales" over "2024", "2025" becomes "Sales 2024
833
             ", "Sales 2025").
834
835
        - Fill in missing headers by inferring from context.
836
        2. Assigning and Validating Headers
837
        Overview:
838
        Once headers are detected, assign them programmatically and ensure they match
839
             expected schema and data types.
840
841
        Implementation:
842
        - Column Naming:
843
        - Set names in code, e.g., df.columns = ["Name", "Date", "Region", "Amount"].
844
        - Schema Mapping:
845
        - Map headers to a standardized schema, using external files or user prompts.
846
847
        - Example:
        - Raw header: "Amt"; Mapped header: "Amount"
848
        - Quality Checks:
849
        - Detect duplicate or empty headers ("Date", "Date" becomes "Date_1", "Date_2").
850
        - Validate each column's expected data type.
851
852
        3. Automation and Usability Enhancements
853
854
        Enhance usability and automation to streamline header extraction and user
855
             interaction.
856
857
858
        Features:
859
        - Freeze Panes:
        - Automatically freeze header rows in Excel for easier navigation.
860
        - Highlighting:
861
        - Use colored formatting to visually distinguish headers.
862
        - Example:
863
        - Yellow fill for header row.
864
        - Documentation:
865
        - Log extraction logic and confidence scores for each detected header.
866
867
        - Build header extraction into ETL pipelines and record process metadata.
868
869
        Block Detection
870
        1. Identifying Block Boundaries
871
        Overview:
872
873
        Block detection segments data into logical units or tables.
874
875
        Methods:
876
        - Boundary Detection:
        - Find empty rows, repeated labels, or formatting changes.
877
878
        | Name | Amount | |-----| | John | 100 | | | <-- Empty row indicates
879
             new block | Name | Amount | | Alice | 200 |
880
        - Machine Learning:
881
882
        - Train classifiers to detect block boundaries based on cell patterns.
883
884
        Advanced:
885
        - Detect nested blocks or hierarchies using indentation or merged cells.
```

- Identify summary blocks with keywords like "Total" or "Summary".

886

```
887
        2. Processing and Tracking Blocks
888
        Overview:
889
        Once blocks are detected, assign IDs and enable block-level analysis.
890
891
        Actions:
892
        - Block ID:
893
        - Assign unique IDs (e.g., Block_001, Block_002).
894
        - Analysis:
895
896
         - Perform group-by or aggregation within each block.
897
        - Example:
        - Sum "Amount" for Block_001: 100 + 150 = 250
898
899
        3. Additional Block Actions
900
        Overview:
901
902
        Enable modular analysis and reporting at the block level.
903
904
        Features:
        - Summary Rows:
905
        - Add computed totals/averages for each block.
906
        - Export/Save:
907
        - Save blocks as separate files or sheets.
908
909
        - Example:
        - Export Block_001 to "block1.csv"
910
911
        Search for Values or Patterns
912
        1. Search Execution Methods
913
        Overview:
914
        Efficiently locate specific values or patterns in your data.
915
916
        Techniques:
917
        - Manual Tools:
918
        - Use Ctrl + F in Excel for quick lookups.
919
        - Programmatic Search:
920
        - Scan all cells using loops or vectorized code.
921
922
        - Example:
923
        - Find all instances of "North" in the "Region" column.
        - Pattern Matching:
924
        - Support exact, wildcard (*Total*), and regex (\d{4}-\d{2}-\d{2} for dates).
925
926
        2. Recording and Highlighting Results
927
        Overview:
928
        Log and visualize search matches for user review.
929
930
        Actions:
931
932
         - Logging:
        - Record coordinates (e.g., Sheet1, Row 3, Col "Region").
933
        - Highlighting:
934
        - Apply conditional formatting to search hits.
935
936
        3. Advanced Search Scenarios
937
938
        Overview:
        Handle complex or large-scale search requirements.
939
940
941
        Scenarios:
942
        - Merged Cells:
        - Search within merged cells or across multiple sheets.
943
944
        - Export:
        - Export found results for further analysis.
945
946
        - Example:
        - Export all rows containing "John" to "john_results.csv"
947
948
949
        Writeback Results
950
        1. Output Placement
        Overview:
951
```

```
Choose where and how to insert results.
 952
 953
         Options:
954
         - Target Columns:
955
 956
         - Select existing or blank columns for output.
         - Appending:
 957
         - Add new columns for flags, counts, or statuses.
958
         - Example:
959
         - Add "Approved_Flag" column next to "Status".
 960
 961
 962
         2. Writing and Styling Results
         Overview:
 963
         Automate and style the output for visibility.
964
965
         Methods:
 966
 967
         - Formulas/Code:
         - Use code (e.g., ws.cell(row, col).value = result) to insert results.
968
969
         - Styling:
         - Bold, borders, or colors for output cells.
 970
         - Example:
971
         - Green fill for "Success", red for "Error".
972
973
         3. Audit and Protection
974
 975
         Overview:
         Maintain the integrity and traceability of results.
976
977
978
         Measures:
         - Lock Columns:
 979
         - Prevent edits to output columns.
 980
         - Timestamps/User Info:
 981
         - Add audit trail for writebacks.
 982
         - Example:
 983
         - "2024-06-01, User: admin"
 984
985
         Difference in State
986
         1. Sheet Comparison
 987
 988
         Overview:
         Identify changes between input and output sheets.
 989
990
991
         Process:
         - Load Sheets:
 992
         - Read both sheets into memory.
 993
         - Compare Cells:
994
         - Detect differences by position and value.
995
 996
 997
         2. Recording and Reporting Differences
         Overview:
998
         Log and report all detected changes.
999
1000
         Actions:
1001
         - Log Mismatches:
1002
         - Record cell coordinates and values.
1003
         - Example:
1004
         - Cell B3: "North" \rightarrow "South"
1005
1006
         - Export Diff Report:
         - List all detected differences for review.
1007
1008
1009
         3. Visualization and Automation
         Overview:
1010
         Make changes visible and automate validation.
1011
1012
         Features:
1013
1014
         - Highlight Changes:
         - Color code changed cells.
1015
         - Automate Checks:
1016
```

```
- Integrate diff comparisons into test scripts.
1017
1018
         Column Selection
1019
         1. Selection Criteria
1020
1021
         Overview:
         Choose relevant columns for analysis.
1022
1023
         Methods:
1024
         - Labels/Indices:
1025
1026
         - Select by name or position.
1027
         - Dynamic Rules:
         - E.g., all numeric columns.
1028
         - Assign Roles:
1029
         - Example: "ID", "Date", "Metric"
1030
1031
1032
         2. Preparation and Validation
         Overview:
1033
         Prepare columns for consistent use.
1034
1035
1036
         Actions:
         - Rename/Relabel:
1037
         - Standardize column names.
1038
1039
         - Validate Types:
         - Ensure columns are of expected type.
1040
         - Example:
1041
         - "Date" column as datetime.
1042
1043
         3. Reusability
1044
1045
         Overview:
         Save and reuse column selections.
1046
1047
         Features:
1048
         - Presets:
1049
         - Save selection profiles.
1050
         - Downstream Use:
1051
         - Use validated columns in subsequent processes.
1052
1053
         Filter Rows
1054
         1. Filtering Methods
1055
         Overview:
1056
         Refine your dataset with filters.
1057
1058
         Techniques:
1059
         - Spreadsheet Tools:
1060
         - Use built-in filters.
1061
1062
         - Code Logic:
         - Filter with code (e.g., df[df['Status'] == 'Approved']).
1063
         - Multiple Criteria:
1064
         - Combine conditions (AND/OR).
1065
         - Example:
1066
         - Status = "Approved" AND Amount > 100
1067
1068
         2. Helper Columns and Complex Filters
1069
1070
         Simplify filtering using helper columns.
1071
1072
         Actions:
1073
         - Helper Columns:
1074
         - Compute intermediate flags.
1075
1076
         - Document Logic:
         - Record filtering rules for audit.
1077
1078
1079
         3. Post-Filter Actions
1080
         Overview:
         Visualize and export filtered data.
1081
```

```
1082
1083
         Features:
         - Highlighting:
1084
         - Grey-out filtered-out rows.
1085
1086
         - Export:
         - Save the filtered dataset.
1087
1088
         Merge Tables
1089
         1. Key-Based Merging
1090
1091
         Overview:
1092
         Combine tables using shared keys.
1093
         Techniques:
1094
         - Join Operations:
1095
         - Use VLOOKUP, JOIN, or code merges.
1096
1097
         - Example:
         - Merge "Customer_ID" from two tables.
1098
         - Align Data:
1099
         - Match on columns like "ID", "Name".
1100
1101
         2. Stack-Based Merging
1102
         Overview:
1103
         Append tables when keys 'arent needed.
1104
1105
         Methods:
1106
         - Vertical Append:
1107
          - Combine rows from similar tables.
1108
         - Deduplicate:
1109
         - Remove duplicate records.
1110
1111
         3. Tracking and Audit
1112
         Overview:
1113
         Track source and unmatched records.
1114
1115
1116
         Actions:
         - Source Column:
1117
1118
         - Add "Source" to indicate origin.
         - Highlight Unmatched:
1119
         - Mark or export mismatched rows.
1120
1121
         Pivot or Unpivot
1122
1123
         1. Pivoting Data
         Overview:
1124
         Summarize data using pivots.
1125
1126
1127
         Methods:
1128
         - PivotTables:
         - Group by row/column dimensions.
1129
         - Example:
1130
         - Sum "Amount" by "Region".
1131
1132
         - Aggregation:
         - Choose SUM, AVG, COUNT, etc.
1133
1134
         2. Unpivoting (Melting) Data
1135
1136
         Overview:
         Reshape data from wide to long format.
1137
1138
1139
         Techniques:
         - Melt Operations:
1140
1141
         - Convert columns into rows.
1142
         - Example:
1143
         | Year | Sales_2019 | Sales_2020 | |-----|
1144
1145
         | Year | Sales_Year | Value |
1146
```

```
- Flexible Restructuring:
1147
1148
          - Selectively unpivot non-ID columns.
1149
          3. Post-Pivot Actions
1150
1151
          Overview:
          Prepare pivoted data for export.
1152
1153
          Features:
1154
          - Flatten Pivot Table:
1155
1156
          - Convert back to flat for further analysis.
          - Reorder/Rename:
1157
          - Clarify pivoted fields.
1158
1159
          Map with Lookup Tables
1160
          1. Mapping Techniques
1161
1162
          Overview:
          Standardize data using lookups.
1163
1164
          Methods:
1165
1166
          - Functions:
          - Use \ensuremath{\text{VLOOKUP}}, merge with dictionaries.
1167
          - Code-to-Label:
1168
          - Example:
1169
          - Code "N" \rightarrowLabel "North"
1170
1171
          2. Application and Fallbacks
1172
1173
          Overview:
          Apply lookups and handle missing values.
1174
1175
          Actions:
1176
          - Apply Mappings:
1177
          - Across selected columns.
1178
1179
          - Handle Missings:
          - Use defaults for missing codes.
1180
1181
          3. Audit and Display
1182
1183
          Overview:
1184
          Ensure mapping transparency.
1185
          Features:
1186
          - Cache Mappings:
1187
1188
          - Store for repeated use.
          - Display Codes/Labels:
1189
          - Show both for clarity.
1190
1191
1192
          Fill Missing Data
          1. Choosing Fill Methods
1193
          Overview:
1194
1195
          Impute missing data appropriately.
1196
1197
          Techniques:
          - Forward/Backward Fill:
1198
          - Fill gaps with prior/next value.
1199
          - Default Values:
1200
          - Use fixed placeholder (e.g., 0, "Unknown").
1201
          - Contextual Example:
1202
          - Dates: Fill missing month with last known month.
1203
1204
          2. Application and Auditing
1205
1206
          Overview:
          Apply fills and flag for review.
1207
1208
1209
          Actions:
          - Targeted Filling:
1210
          - Apply to specific columns/rows.
1211
```

```
- Flag Filled Cells:
1212
1213
          - Highlight for later review.
1214
          3. Documentation
1215
1216
          Overview:
1217
          Keep fill logic transparent.
1218
          Features:
1219
          - Record Logic:
1220
1221
          - Document assumptions and methods.
1222
          - Audit Trail:
          - Track all changes.
1223
1224
          Flag Rows or Cells
1225
          1. Defining Flag Rules
1226
1227
          Overview:
          Establish criteria for flagging.
1228
1229
          Examples:
1230
1231
          - Simple Rule:
          - Flag where Amount < 0
1232
          - Complex Rule:
1233
          - Flag where Status = "Pending" and Amount > 1000
1234
1235
          2. Applying Flags
1236
          Overview:
1237
          Insert flags and summarize.
1238
1239
          Actions:
1240
          - Flag Column:
1241
1242
          - Add "Flag" column with "Yes"/"No".
          - Export Flagged Rows:
1243
1244
          - Save for further inspection.
1245
          3. Advanced Flagging
1246
          Overview:
1247
1248
          Use multiple criteria and document.
1249
          Features:
1250
1251
          - Multi-Criteria:
          - Combine several rules for granular checks.
1252
1253
          - Notes:
          - Document flagging rationale.
1254
1255
1256
          Sort Data
1257
          1. Setting Sort Criteria
1258
          Overview:
          Organize data for analysis.
1259
1260
          Options:
1261
1262
          - Sort Columns:
          - By value, ascending/descending.
1263
          - Multi-Level:
1264
          - E.g., sort by "Region", then by "Amount".
1265
1266
          2. Applying Sorts
1267
          Overview:
1268
          Implement sorting programmatically or manually.
1269
1270
1271
          Methods:
          - Spreadsheet Tools:
1272
          - Built-in sort features.
1273
1274
          - Code:
          - E.g., df.sort_values(['Region', 'Amount'])
1275
1276
```

```
3. Post-Sort Actions
1277
1278
          Overview:
          Finalize sorted data.
1279
1280
1281
          Actions:
          - Renumber Rows:
1282
          - Update indices.
1283
          - Highlight Extremes:
1284
          - Mark top/bottom values.
1285
1286
1287
          Validate Data
          1. Validation Checks
1288
          Overview:
1289
          Ensure data meets required standards.
1290
1291
1292
          Checks:
          - Type:
1293
1294
          - Ensure numeric columns contain numbers.
1295
          - Range:
          - E.g., "Amount" > 0.
1296
          - Pattern:
1297
          - Date columns match YYYY-MM-DD.
1298
          - Business Rule Example:
1299
          - "Start Date" < "End Date"
1300
1301
          2. Marking and Reporting
1302
1303
          Overview:
          Visualize and report errors.
1304
1305
          Actions:
1306
1307
          - Highlight Invalids:
          - Color-code errors.
1308
1309
          - Export Summary:
          - Table of error counts and locations.
1310
1311
          3. Integration in Workflow
1312
1313
          Make validation a routine part of processing.
1314
1315
          Features:
1316
          - Pre-Processing Step:
1317
1318
          - Validate before analysis.
          - Automation:
1319
          - Integrate into data pipelines.
1320
1321
1322
          Split Sheets or Data
          1. Defining Split Rules
1323
          Overview:
1324
1325
          Segment data for modular analysis.
1326
1327
          Methods:
          - By Category:
1328
          - E.g., split by "Region".
1329
          By Date Range:E.g., split by year.
1330
1331
1332
          2. Exporting Segments
1333
1334
          Overview:
          Save segments for separate use.
1335
1336
1337
          Actions:
          - Export Files:
1338
          - "North_Region.csv", "South_Region.csv"
1339
          - Consistent Formatting:
1340
          - Ensure identical columns and styling.
1341
```

```
1342
          3. Automation and Documentation
1343
1344
          Overview:
          Automate splitting and track provenance.
1345
1346
          Features:
1347
          - Automation:
1348
          - Use scripts/macros for repeated splits.
1349
           Documentation:
1350
           Record rules and export logs.
1351
```

D Qualitative Analysis of BREW-Generated Knowledge Bases

This section presents a comprehensive qualitative analysis of knowledge bases generated through the BREW technique applied to two distinct agent training environments: OSWorld and τ^2 Bench described in the section before. The analysis examines knowledge representation patterns, procedural sophistication, and domain-specific learning characteristics extracted from CUA agent behaviors, providing insights into the effectiveness and scope of knowledge distillation techniques across diverse task environments.

D.1 Cross-Domain Knowledge Base Analysis

D.1.1 Base Structure & Organization

1353

1360

1361

1389

Schema Consistency and Evolution: Both knowledge bases demonstrate consistent structural schemas, though adapted to their respective domains. The OSWorld KB employs a four-part schema (contextual triggers, procedural steps, extended capabilities, concrete instantiation), while the τ^2 Bench KB extends this to a five-part structure, adding explicit purpose rationale ("Why to use it"). This evolution suggests that BREW adapts its extraction patterns to domain-specific requirements—conversational commerce demands explicit justification for actions due to customer interaction contexts.

Taxonomic Organization Principles: The OSWorld KB reveals a capability-based taxonomy organized around computational tasks: file operations, document processing, inter-application workflows, and data visualization. Each category represents a distinct computational domain with specific tool requirements and interaction patterns. In contrast, the τ^2 Bench KB employs a **lifecycle-based taxonomy** structured around transactional states: order creation, modification, fulfillment, and post-delivery operations. This organizational difference reflects fundamental domain characteristics—desktop automation focuses on tool orchestration, while conversational commerce centers on process management.

Hierarchical Task Decomposition: Both KBs demonstrate sophisticated hierarchical reasoning, but through different decomposition strategies. OSWorld exhibits **technical decomposition**, breaking complex operations like "Create Charts from Data" into constituent technical steps (data selection, chart insertion, customization, formatting). τ^2 Bench shows **process decomposition**, structuring operations like order modification into authentication, validation, confirmation, and execution phases. This suggests BREW successfully identifies domain-appropriate decomposition strategies rather than applying uniform patterns.

Knowledge Boundary Definition: Both KBs explicitly encode operational boundaries, but through contrasting mechanisms. OSWorld boundaries are **capability-constrained**—determined by available applications and system resources. τ^2 Bench boundaries are **policy-constrained**—explicitly defined through "Deny Unsupported Request" patterns and escalation protocols. This difference highlights how knowledge extraction adapts to domain-specific constraint types.

D.1.2 Procedural Knowledge Grounding

Context-Dependent Action Selection: Both domains demonstrate sophisticated context awareness, but grounded in different environmental factors. OSWorld exhibits application-context sensitivity, where identical operations (e.g., image insertion) require different procedures across LibreOffice Writer, Impress, GIMP, and Thunderbird. The agent learned application-specific affordances and

interaction patterns rather than generic command sequences. τ^2 Bench demonstrates **state-context** sensitivity, where available actions depend on order status (pending vs. delivered), payment methods, and authentication levels. This reveals learned understanding of business process constraints and temporal operation windows.

Error Prevention and Validation Workflows: Both KBs incorporate sophisticated error prevention mechanisms, but grounded in domain-specific failure modes. OSWorld emphasizes technical validation: file integrity checks ("confirm the exported file opens correctly"), application state verification, and multi-step confirmation for irreversible operations. τ^2 Bench emphasizes transactional validation: authentication cascades, confirmation dialogues with standardized templates, and explicit user consent protocols. The emergence of defensive programming practices across both domains suggests these represent fundamental principles of reliable agent behavior.

State-Dependent Decision Logic: The procedural knowledge in both domains demonstrates sophisticated state machine reasoning. OSWorld exhibits application state awareness—understanding when applications are ready for input, when files are loaded, and when operations can be safely executed. Window management and application switching reveal learned understanding of desktop metaphors and resource constraints. τ^2 Bench demonstrates business process state awareness—finite state machine reasoning where order lifecycle states determine available operations. The agent learned that pending orders enable modification while delivered orders unlock return workflows, indicating internalized understanding of business logic constraints.

Security and Authentication Grounding: While OSWorld operates in a trusted desktop environment with minimal explicit security concerns, τ^2 Bench reveals pervasive authentication-first paradigms. Nearly every transactional operation begins with identity verification through email, name, and zip code combinations. The KB demonstrates graduated security reasoning: information retrieval requires basic authentication while financial transactions trigger rigorous verification protocols. This contrast highlights how procedural knowledge adapts to domain-specific security requirements.

Cross-Application vs. Cross-Process Orchestration: OSWorld demonstrates technical orchestration—coordinating multiple applications (Chrome, LibreOffice suite, File Manager, GIMP) to accomplish complex workflows. The "Navigate Between Applications" section reveals learned behaviors for window management, application switching, and resource coordination. τ^2 Bench exhibits process orchestration—coordinating authentication, validation, confirmation, and execution phases across different operational contexts. Both forms of orchestration require sophisticated temporal reasoning and constraint management, but applied to different environmental complexity types.

Failure Mode Internalization: Both KBs reveal learned understanding of domain-specific failure modes. OSWorld incorporates file validation, application crash recovery suggestions, and verification steps for critical operations. τ^2 Bench includes explicit escalation protocols ("Transfer to Human Agent"), policy compliance mechanisms, and irreversibility warnings for financial operations. The consistent emergence of failure-aware procedures suggests that agents successfully internalize risk assessment and mitigation strategies during training.

1426

1427

1428

1429

1430

1431

Domain-Specific Communication Patterns: The procedural knowledge reveals distinct communication paradigms appropriate to each domain. OSWorld procedures are **task-oriented** with minimal user interaction—focusing on efficient command execution and verification. τ^2 Bench procedures are **dialogue-oriented** with standardized customer interaction templates, confirmation protocols, and expectation management communications. This adaptation demonstrates that BREW extracts not just procedural logic but domain-appropriate interaction modalities.

The cross-domain analysis reveals that BREW successfully extracts procedural knowledge that is both structurally consistent (following learnable organizational patterns) and contextually grounded (adapted to domain-specific constraints, failure modes, and interaction requirements). This dual capability suggests significant potential for knowledge transfer across related domains while maintaining appropriate domain-specific adaptations.