ACON: OPTIMIZING CONTEXT COMPRESSION FOR LONG-HORIZON LLM AGENTS

Anonymous authors

000

001

002003004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

023

024

025

026

027

028

031

033

034

035

037

040

041

042

043

044

045

047

048

051

052

Paper under double-blind review

ABSTRACT

Large language models (LLMs) are increasingly deployed as agents in dynamic, real-world environments, where success requires both reasoning and effective tool use. A central challenge for agentic tasks is the growing context length, as agents must accumulate long histories of actions and observations. This expansion raises memory costs and reduces token efficiency in long-horizon tasks, yet prior work on context compression has mostly focused on single-step tasks or narrow applications. We introduce **Agent Context Optimization** (ACON), a unified framework that optimally compresses both environment observations and interaction histories into concise yet informative condensations. ACON leverages compression guideline optimization in natural language space: given paired trajectories where full context succeeds but compressed context fails, capable LLMs analyze the causes of failure, and the compression guideline is updated accordingly. Furthermore, we propose distilling the optimized LLM compressor into smaller models to reduce the overhead of the additional module. Experiments on AppWorld, OfficeBench, and Multi-objective QA show that ACON reduces memory usage by 26-54% (peak tokens) while largely preserving task performance, preserves over 95% of accuracy when distilled into smaller compressors, and enhances smaller LMs as long-horizon agents with up to 46% performance improvement.

1 Introduction

Large language models (LLMs) have become the backbone of AI agents across diverse real-world tasks, leveraging language knowledge and reasoning to plan, act, and adapt in dynamic environments (Yao et al., 2023). These tasks often unfold over extended horizons: the agent must gather information, invoke tools, and revise plans based on feedback. In such settings, context is not auxiliary but foundational, as tool APIs and data formats lack standardization (Anthropic, 2024a), and even personalized environments differ in file structures or account data. Agents must therefore accumulate and maintain records of prior actions, observations, and world state. Losing this information, such as an email identifier, file version, or API format, can derail task success.

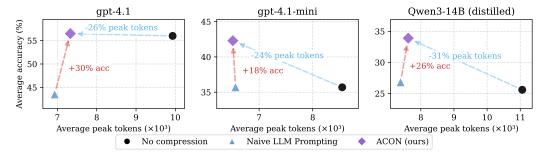


Figure 1: Accuracy–Peak tokens trade-off on AppWorld (Trivedi et al., 2024). We compare average accuracy versus peak input tokens in history compression. ACON (ours) reduces cost while preserving accuracy for the large model (gpt-4.1) relative to a naive prompting baseline, and even improves accuracy on smaller models (gpt-4.1-mini and Qwen-14B). More results are in Section 4.

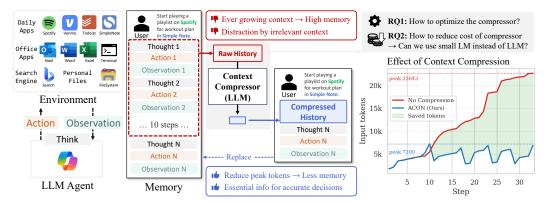


Figure 2: **Motivation: Unbounded context in LLM agents.** As LLM agents interact with environments, actions and observations continuously accumulate, leading to ever-growing contexts that incur high memory usage as in the red line on the right plot. This motivates *Agent Context Optimization* (ACON), which optimally compresses histories and observations into concise summaries, reducing peak tokens and memory as in the blue line on the right plot.

The centrality of context, however, makes it a bottleneck. As interactions accumulate, contexts grow unbounded as in Figure 2, creating two major challenges. First, the inference memory cost of transformers scales with the context length, which becomes prohibitive in long-horizon tasks. Second, excessively long contexts dilute relevant information, distracting the model with outdated or extraneous details (Shi et al., 2023). These issues make effective context management and compression indispensable. Existing compression approaches only partially address this need. Dialogue-oriented systems rely on session-level summarization or tiered memories suitable for conversational coherence but inadequate for multi-step carry-over (Packer et al., 2023). Document-centric methods in long-context QA or in-context learning (Li et al., 2023; Jiang et al., 2024) assume single-step reasoning where context can be discarded after producing an answer. While effective in their domains, these strategies fall short for complex multi-step tasks where success depends on retaining structured signals across many interactions.

The challenge becomes particularly critical in productivity scenarios (e.g., email management, document processing, or workflow automation) where agents must coordinate across heterogeneous tools and maintain precise state information. Unlike simpler agent tasks, these environments demand preservation of diverse signal types: factual history, action—outcome relationships, evolving environment states, success preconditions, and future decision cues. Naive strategies like token truncation or generic summarization easily lose critical details essential for long-horizon reasoning. Recent agent-focused compression methods (Deng et al., 2023; Lee et al., 2025; Smith, 2025) either specialize narrowly or rely on brittle heuristics, limiting their applicability across the full spectrum of multi-step agent tasks.

We address these challenges with **Agent Context Optimization** (**ACON**), a unified framework for systematic and adaptive context compression. Our study yields three key findings. First, task- and environment-specific guidelines enable consistent context compression across diverse agents without sacrificing performance. Second, optimized contexts not only reduce memory costs but also improve decision quality, allowing smaller LLMs to act more effectively. Third, high-quality compressors can be distilled into even smaller models, reducing overhead and improving deployability.

ACON dynamically condenses environment observations and interaction histories into concise, informative representations. Rather than handcrafting prompts, we introduce a guideline optimization pipeline that refines compressor prompts via failure analysis in natural language space (Pryzant et al., 2023; Khattab et al., 2024; Yüksekgönül et al., 2025; Han et al., 2025), ensuring that critical environment-specific and task-relevant information is retained. Importantly, ACON is gradient-free, requiring no parameter updates, making it directly usable with closed-source or production models. We further distill optimized compressors into smaller models for cost-efficient deployment.

We validate ACON on three multi-step agent benchmarks: AppWorld (Trivedi et al., 2024), OfficeBench (Wang et al., 2024b), and Multi-objective QA (Kwiatkowski et al., 2019; Zhou et al., 2025), each requiring 15+ interaction steps. Our empirical results demonstrate clear advantages of ACON: (i) lowers memory usage by 26–54% (peak tokens) while largely maintaining task perfor-

mance; (ii) enables effective distillation of the context compressor into smaller models, preserving 95% of the teacher's accuracy across all benchmarks, thereby reducing the overhead of the additional module; (iii) allows small LMs to function more effectively as agents, improving performance by 32% on AppWorld, 20% on OfficeBench, and 46% on Multi-objective QA by mitigating the distraction of long contexts through ACON. Our result highlights on AppWorld benchmark are in Figure 1.

To summarize, our work makes the following contributions:

- We propose **Agent Context Optimization** (ACON), a framework for compressing both environment observations and interaction histories, tailored to multi-step, long-horizon agentic tasks.
- We develop a failure-driven, task-aware compression guideline optimization. Our approach is entirely **gradient-free** and readily applicable to any LLMs, including API-based models.
- We provide a cost-efficient solution by distilling optimized compressors into smaller models, preserving over 95% of the teacher's performance while reducing module overhead.
- We validate ACON on AppWorld, OfficeBench, and Multi-objective QA, showing that it reduces memory usage by 26-54% (peak tokens) while preserving task success with large LLMs, and enabling small LMs to achieve 20-46% performance improvements.

2 RELATED WORKS

Long-horizon LLM agents. Large language model (LLM) agents extend pretrained models beyond static single-step reasoning tasks (e.g., RAG-based QA, math problem solving, or code generation) to interactive decision-making in dynamic environments (Yao et al., 2023; Wang et al., 2024a; Team et al., 2025; OpenAI, 2025a). Unlike chatbots or solvers that return an answer in one pass, agents must iteratively observe their surroundings, select tools, and execute actions while revising their plans based on feedback (Shridhar et al., 2021; Jimenez et al., 2024; Zhou et al., 2024; Wei et al., 2025; Xie et al., 2024; Bonatti et al., 2024). Recent work highlights the importance of *long-horizon LLM agents*, which tackle tasks that unfold over dozens to hundreds of steps and require coordination across multiple applications and tools (Kwa et al., 2025; Trivedi et al., 2024; Wang et al., 2024b; 2025b). A central challenge in these scenarios lies in managing the *dynamic long context*, where the agent must retain multi-step interaction histories and handle diverse observations produced by heterogeneous environments.

Context compression for LLMs. Managing this ever-growing context has been a longstanding challenge, and a variety of approaches have been proposed to compress LLM inputs. Prior works on context compression can be broadly grouped into three directions: document- or retrieval-based compression (Seo et al., 2025; Li et al., 2023; Xu et al., 2024; Yoon et al., 2024; Zhou et al., 2025; Jiang et al., 2024; Shandilya et al., 2025), dialogue memory summarization (Xu et al., 2025; Maharana et al., 2024; Wang et al., 2025a), and low-level KV cache compression (Zhang et al., 2025). While each line of research has demonstrated benefits in its respective setting, they remain insufficient for the dynamic and heterogeneous contexts required by long-horizon agents.

Beyond these directions, a few recent studies have attempted to compress context specifically for LLM agents (Deng et al., 2023; Lee et al., 2025; Smith, 2025; Yang et al., 2024b). However, these approaches either rely on naive prompting or target narrow domains such as the accessibility tree in web browsing, limiting their generality. In contrast, our work introduces a universal *agent context optimization* framework that is applicable to arbitrary agents. It supports both history and observation compression, and incorporates a generalizable optimization methodology that is agnostic to the underlying model, making it applicable to both open-source and proprietary API-based LLMs.

A closely related line of work is MEM1 (Zhou et al., 2025), which learns a compression policy jointly with the agent via reinforcement learning. This couples reasoning and compression and restricts applicability to models with trainable weights, whereas our method can work without weight updates. A more detailed comparison is provided in Appendix C.

3 AGENT CONTEXT OPTIMIZATION (ACON)

We present Agent Context Optimization (ACON), a unified framework for optimized history and observation compression in long-horizon LLM agents. We begin by formulating the agent task and

defining context cost in Section 3.1. Next, in Section 3.2, we introduce generative compression with LLMs for both history and observation, and formalize the associated optimization objective and its challenges. We then propose our optimization method in Section 3.3, followed by a distillation that enables smaller models for compressions to improve cost efficiency (Section 3.4).

3.1 PROBLEM FORMULATION

Task. An agentic task is formulated as a Partially Observable Markov Decision Process (POMDP) $\mathcal{E} = \langle \mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{T}, \mathcal{R} \rangle$ with state space \mathcal{S} , action space \mathcal{A} , and observation space \mathcal{O} . The transition function $\mathcal{T}(s,a) \to (s',o)$ is deterministic, and it is determined by the implementation of the environment. Specifically, it executes an action $a \in \mathcal{A}$ in the environment and returns the next state and observation. The reward function \mathcal{R} returns the reward given the terminal state s_T . The terminal state is arrived when the transition function receives the special action (e.g., finish_task).

An LLM agent interacts with the environment to get information for making a decision to achieve a given task o_0 through multiple steps. For each step t, the LLM \mathcal{M} generates the action a_t followed by its reasoning at each step (Yao et al., 2023; Wang et al., 2024a) given the interaction history $\mathbf{h}_{t-1} = (o_0, a_0, o_1, a_1, \cdots, o_{t-1}, a_{t-1})$ and the latest observation o_t :

$$\mathcal{M}(o_t, \boldsymbol{h}_{t-1}; \theta, \mathcal{P}_{\mathsf{agent}}) \mapsto a_t,$$
 (1)

where θ refers to the pre-trained parameters of the LLM and \mathcal{P}_{agent} is the prompt that consists of a general environment description, tools, output format, and few-shot examples in natural language.

Cost function for context. We assume that the LLM agent's parameters θ and the task and system prompt $\mathcal{P}_{\text{agent}}$ are fixed. We define a cost function \mathcal{C} that measures the cost of encoding the dynamic context during action generation at each step such as $\mathcal{O}(n)$ computational cost of a transformer for decoding given n input tokens. The cost function takes the interaction history h_{t-1} , and the latest observation o_t as input and returns the per-step cost:

$$C(\boldsymbol{H}) = \sum_{t=1}^{T} C(\boldsymbol{h}_{t-1}, o_t),$$
 (2)

where C is the total cost of completing the task, $\mathbf{H} = \{\mathbf{h}_{t-1}, o_t\}_{t=1}^T$ denotes the sequence of history and observation of each step. Typically, C is proportional to the summation of token lengths of action and observations in each step, \mathbf{h}_{t-1} and o_t . While the prompt cost is static and can be budgeted in advance, the costs from interaction histories are **unbounded**, leaving the user with only two options: terminate the task early or truncate the context heuristically to a maximum length. This raises the central question: how can we compress context more effectively than such heuristics?

3.2 HISTORY & OBSERVATION COMPRESSION WITH LLMS

To address this challenge, we use an LLM $f(\cdot; \phi, \mathcal{P})$, parameterized by pre-trained weights ϕ and a compression guideline \mathcal{P} , to minimize context cost defined in Equation 2 (e.g., summarize the given interaction history). As in Equation 1, the LLM receives two inputs at each step: the interaction history h_{t-1} and the latest observation o_t . This introduces two options for context compression:

History compression. The interaction history accumulates both environment observations and agent actions. In long-horizon tasks, this history can grow excessively large. To manage its length, we apply history compression only when the history length exceeds a predefined threshold T_{hist} :

$$\mathbf{h}'_t = f(\mathbf{h}_t; \phi, \mathcal{P}_{\mathsf{hist}}) \text{ if } |\mathbf{h}_t| > T_{\mathsf{hist}}, \quad \mathbf{h}_t \text{ otherwise.}$$
 (3)

The compressed history h'_t replaces the raw history in Equation 1. This selective compression ensures that the overhead of invoking the compressor is incurred only when necessary (Smith, 2025).

Latest observation compression. Given an action a, the environment returns an observation o according to the transition function $\mathcal{T}(s,a) \to (s',o)$. We similarly apply observation compression only when the observation length exceeds a threshold T_{obs} :

$$o'_t = f(o_t, \mathbf{h}_{t-1}; \phi, \mathcal{P}_{\text{obs}}) \text{ if } |o_t| > T_{\text{obs}}, \quad o_t \text{ otherwise.}$$
 (4)

Figure 3: **Compression Guideline Optimization.** Feedback is generated by contrasting successful trajectories (no compression) with failed ones (with compression). The collected feedback is then used by LLM to refine the compression guidelines.

This mechanism avoids unnecessary overhead when o_t is already short, while still reducing redundant or distracting content in long observations (Xu et al., 2024; Deng et al., 2023; Lee et al., 2025). The compressed one o'_t replaces the raw one in Equation 1 and is stored in the interaction history h.

In both cases, the compressor LLM selects information to preserve based on its learned prior knowledge of importance. However, there is **no guarantee** that the salient details required for successful task completion are retained. The agent context effectively serves as a **world model of the environment**, encompassing diverse forms of information such as causal relations (*e.g.*, email leaves drafts), evolving states (*e.g.*, account balance), preconditions (*e.g.*, login required), and task-relevant decision cues (*e.g.*, due dates). Effective context compression must therefore accommodate this heterogeneous and dynamic nature of agent context, ensuring that the most critical signals are preserved for long-horizon reasoning and task success.

Optimization objective. We optimize the compressor parameters $\psi \triangleq (\phi, \mathcal{P})$ to maximize task reward while minimizing context cost. At each step t, the compressor produces either a compressed history $h'_t = f_{\mathsf{hist}}(h_t; \psi)$ or observation $o'_t = f_{\mathsf{obs}}(o_t, h_{t-1}; \psi)$. Let the compressed context be

$$H'(\psi) = \{h'_{t-1}, o'_t\}_{t=1}^T, \qquad C(H'(\psi)) = \sum_{t=1}^T C(h'_{t-1}, o'_t).$$
 (5)

With the agent $\mathcal{M}(\cdot; \theta, \mathcal{P}_{\text{agent}})$ fixed, the environment induces a trajectory $\tau(\psi)$ and terminal state $s_T(\psi)$ when the agent conditions on $H'(\psi)$. Our learning objective is

$$\max_{\psi} \underbrace{\mathbb{E}\left[\mathcal{R}(s_T(\psi))\right]}_{\text{maximize}} - \lambda \underbrace{\mathbb{E}\left[C(\mathbf{H}'(\psi))\right]}_{\text{minimize}}, \qquad \lambda \ge 0, \tag{6}$$

where λ is a multiplier and the expectations are over tasks.

Challenges. The optimization objective in Equation 6 is difficult to optimize in practice because there is no gold supervision for compression, the reward is sparse and only revealed at the end of the trajectory, and the context cost is defined over discrete quantities, which precludes direct gradient computation. While these properties naturally motivate reinforcement learning (RL) (Sutton & Barto, 2018), applying RL to this setting introduces additional obstacles: (1) updating the parameters ϕ of a LLM with RL can be computationally prohibitive, (2) environment roll-outs are extremely expensive since each reward requires multi-step executions of both agent and compressor LLMs, and (3) policy gradient estimates suffer from high variance because compression quality is only indirectly evaluated through eventual task success.

3.3 OPTIMIZING COMPRESSION GUIDELINES

To overcome these challenges, we propose to optimize **compression guidelines** \mathcal{P} (natural language prompts) for context compression, rather than fine-tuning model parameters ϕ . Trajectories under compressed contexts provide *dense signals* about the quality of compression. For example, if the agent fails with compressed context while succeeding without compression, this indicates that the compressed context may have lost crucial information. Such trajectory-level comparisons yield richer feedback than scalar rewards (e.g., binary task success).

We instantiate this idea as prompt optimization using an LLM as the optimizer, where the natural language prompt \mathcal{P} is refined via feedback expressed in natural language (Yang et al., 2024a; Yüksekgönül et al., 2025; Khattab et al., 2024). We introduce **compression guideline optimization** based on *contrastive task feedback*.

On the training set $\mathcal{D}_{\text{train}}$, we run the LLM agent both without and with context compression to obtain baseline context H and compressed context H'. We collect tasks where the agent succeeds with H but fails with H', forming a contrastive subset $\mathcal{D}_{\text{cont}}$. For each task in $\mathcal{D}_{\text{cont}}$, we query a optimizer LLM with the context before and after compression to obtain natural language feedback:

$$Feedback = LLM(Feedback Instruction, H, H').$$
 (7)

This feedback serves as a natural language gradient (Yüksekgönül et al., 2025), indicating how the compression guideline \mathcal{P} should be refined. We then aggregate feedback from multiple trajectories:

$$\mathcal{P}^{(1)} = \text{LLM}(\text{Update Instruction}, \mathcal{P}^{(0)}, \|_{i=1}^n \text{Feedback}_i), \tag{8}$$

where \parallel is concatenation of feedbacks from each task, which corresponds to a batch optimization step in textual gradient descent (Yüksekgönül et al., 2025). We also generate multiple candidate prompts $\{\mathcal{P}_k^{(1)}\}$, evaluate them on \mathcal{D}_{cont} , and select the best-performing one. We refer this process as *utility maximization* step $\underline{\mathtt{UT}}$ as it primarily maximizes the first term (task reward) of Equation 6.

However, optimizing only for reward may neglect the context cost (second term in Equation 6). To address this, motivated by alternating optimization, we perform a second iteration that conditions only on successful task with compressed context, asking the LLM to generate feedback about which information was actually used during execution. This refines $\mathcal{P}^{(1)} \to \mathcal{P}^{(2)}$, encouraging shorter yet sufficient contexts. We refer this additional process as *compression maximization* step $\overline{\text{CO}}$ as it minimizes the second term (context cost) of Equation 6.

We illustrate overall process in Figure 3. Algorithm 1 and prompts are in Appendix B.

3.4 DISTILLING CONTEXT COMPRESSION INTO SMALL MODELS

While compression guideline optimization enables effective compression, repeatedly invoking the large LLM for compression adds substantial overhead. To reduce this cost, we **distill the compressor into a smaller model**. The teacher with optimized guideline \mathcal{P}^* (parameters ϕ_T) generates compressed outputs y from input x, which supervise the student (parameters ϕ_S). We train the student with a cross-entropy objective (Kim & Rush, 2016) with input-output pair (x, y), where $(x, y) = (h_t, h'_t)$ for Equation 3 or $(x, y) = ((h_{t-1}, o_t), o'_t)$ for Equation 4:

$$\min_{\phi_{\mathsf{S}}} \mathbb{E}_{(\boldsymbol{x},\boldsymbol{y}) \sim \mathcal{D}_{\mathsf{train}}^{+}} \left[-\sum_{n=1}^{L_{\boldsymbol{y}}} \log f(\boldsymbol{y}_{n} \mid \boldsymbol{x}, \boldsymbol{y}_{< n}; \phi_{\mathsf{S}}, \mathcal{P}^{*}) \right], \tag{9}$$

where \mathcal{D}_{train}^+ denotes tasks where the teacher succeeds with compressed context.

Once trained, the student replaces the teacher during inference, decoupling decision making from compression. This two-stage pipeline, guideline optimization then distillation, achieves effective compression with a much smaller model ($|\phi_T| \gg |\phi_S|$):

$$f(\cdot; \phi_{\mathsf{T}}, \mathcal{P}) \xrightarrow{\text{prompt optimization}} f(\cdot; \phi_{\mathsf{T}}, \mathcal{P}^*) \xrightarrow{\text{distillation}} f(\cdot; \phi_{\mathsf{S}}, \mathcal{P}^*). \tag{10}$$

4 EXPERIMENTS

We evaluate ACON on three challenging benchmarks that require multi-step interactions across diverse domains. Our experiments are designed to address the following key questions:

- How well does ACON improve token efficiency while preserving performance? (Section 4.2)
- Does distilling the compressor reduce its size while maintaining agent performance? (Section 4.3)
- Can ACON help small, distilled LM agents perform better under long contexts? (Section 4.4)

Table 1: Results across different difficulty levels on **Appworld** benchmark (test-normal). Each block reports accuracy (task goal completion score), steps, peak input tokens (10^3) , and dependency (10^6) for agents. Best results in each column are highlighted in bold. Rows in blue background indicate the results from **ours**. ACON consistently improves accuracy while reducing peak tokens and dependency, with ACON <u>UTCO</u> achieving the best overall performance.

Method		Average	(168)			Easy (57)			edium (48	3)	Hard (63)		
cuiou	Acc. ↑	Steps ↓	Peak ↓	Dep.↓	Acc. ↑	Peak ↓	Dep.↓	Acc. ↑	Peak ↓	Dep.↓	Acc. ↑	Peak ↓	Dep.↓
				Agent: 9	rpt-4.1	/ Compre	essor: gp	t-4.1					
No compression	56.0	16.14	9.93	5.96	80.7	7.57	2.98	47.9	10.10	5.36	39.7	11.95	9.11
History Compres	ssion												
FIFO	45.8	28.48	6.73	5.69	84.2	5.85	2.89	39.6	7.26	6.24	15.9	7.14	7.80
Retrieval	27.4	33.17	8.39	6.68	61.4	7.40	3.97	12.5	8.74	7.72	7.9	9.02	8.33
LLMLingua	39.3	24.42	7.50	6.37	66.7	6.38	3.04	37.5	8.04	7.39	15.9	8.09	8.59
Prompting	43.5	24.01	6.93	5.29	66.7	6.36	2.84	41.7	7.10	5.36	23.8	7.31	7.48
ACON UT	51.2	20.92	7.17	4.49	77.2	6.45	2.43	50.0	7.39	4.47	28.6	7.65	6.37
ACON UTCO	56.5	22.82	7.33	4.69	86.0	7.09	2.84	56.2	7.48	4.43	30.2	7.44	6.55
Observation Con	npression	1											
LLMLingua	32.1	18.16	8.17	6.01	54.4	5.78	2.33	29.2	8.24	5.23	14.3	10.29	9.92
Prompting	42.3	17.38	6.58	4.09	64.9	4.92	1.88	35.4	6.96	4.11	27.0	7.79	6.07
ACON UT	47.0	16.67	7.62	5.08	70.2	5.87	2.21	45.8	7.79	5.00	27.0	9.07	7.73
ACON UTCO	53.6	18.12	7.43	4.93	82.5	5.66	2.63	47.9	7.30	4.43	31.8	9.14	7.50

4.1 EXPERIMENTAL SETUP

Benchmarks & Metrics. We focus on long-horizon agentic task benchmarks that require 10+ interaction steps on average: (1) **AppWorld (Trivedi et al., 2024)**: Main benchmark with 9 simulated apps (e.g., Venmo, Spotify, SimpleNote) and ~ 100 simulated users. Performance is measured by task completion score. (2) **OfficeBench (Wang et al., 2024b)**: Productivity tasks across 6 apps (e.g., Word, Excel, Email), operating on simulated documents. Performance is measured by benchmark-defined accuracy functions. (3) **8-objective QA (Kwiatkowski et al., 2019; Zhou et al., 2025**): QA benchmark where agents interact with a search tool to answer 8 questions and output a consolidated answer set. Performance is the average of Exact Match (EM) and F1 scores across 8 questions.

In addition to task-specific performance metrics, we report three token efficiency metrics following prior work (Zhang et al., 2025; Zhou et al., 2025): (1) Steps: The average number of interaction steps per task. (2) Peak Tokens: The maximum context length encountered across all steps. (3) Dependency: The cumulative dependency of each generated action on prior tokens, measuring how much generation relies on the context history. Full details are provided in the Appendix B.

Baselines. (1) No Compression: full uncompressed context. (2) FIFO: keep the most recent k interactions, discarding earlier ones (Yang et al., 2024b). (3) **Retrieval:** select k past interactions most similar to the current query via embedding search (Xu et al., 2025). (4) **LLMLingua:** extractive compression with an encoder-only LM (Jiang et al., 2023; Pan et al., 2024). (5) **Prompting:** naive baseline using a general compression instruction (Smith, 2025; Lee et al., 2025).

Our Methods. We evaluate two versions of ACON. (1) ACON <u>UT</u> utilizes an *optimized guideline* for context compression after utility maximization step. (2) ACON <u>UTCO</u> applies compression maximization step. CO after utility maximization <u>UT</u>, aiming for shorter but informative compression.

4.2 Overall Performance and token efficiency on LLMs

In Table 1 and Table 2, we first evaluate ACON using on gpt-4.1 (OpenAI, 2025b) for both agent and compressor, which already achieves strong results on three long-horizon benchmarks.

For history compression, as shown in Table 1, on AppWorld, ACON reduces peak tokens by over 25% while preserving the accuracy of the no compression upper bound, outperforming all baselines that suffer severe degradation on medium and hard tasks spanning longer steps. On OfficeBench (Table 2a), ACON lowers peak context size by nearly 30% while maintaining accuracy above 74%. On 8-objective QA (Table 2b), ACON even surpasses the no compression baseline in EM/F1 while reducing peak tokens and dependency by 54.5% and 61.5%, respectively. For observation compression, ACON consistently outperforms all baselines confirming that compression guideline optimization is effective for compressing not only history but also raw observations.

Table 2: Results on **OfficeBench** and **8-objective QA** benchmarks. We report performance metrics (acc/EM/F1) along with steps, peak input tokens (10^3) , and dependency (10^6) . Best values are in **bold**. Rows in blue are **ours**. ACON consistently improves accuracy/efficiency trade-offs.

	(a) Offic	eBench					(b) 8	-objecti	ive QA		
Method	Acc. ↑	Steps ↓	Peak ↓	Dep. ↓		Method	ЕМ↑	F1 ↑	Steps ↓	Peak ↓	Dep. ↓
Agent: gpt	-4.1/C	ompresso	r: gpt-4	.1		Agent:	gpt-4.	1 / Com	pressor: g	pt-4.1	
No Compression	76.84	11.52	7.27	4.43		No compression	0.366	0.488	15.78	10.35	3.32
History Compression						History Compres	ssion				
FIFO	67.37	12.26	4.02	2.64		FIFO	0.293	0.388	19.26	5.09	2.51
Retrieval	65.26	16.20	4.33	2.06		Retrieval	0.331	0.438	20.06	5.11	2.62
LLMLingua	70.53	10.89	4.65	1.85		LLMLingua	0.363	0.481	17.68	5.68	2.24
Prompting 71.58 10.13 4.40 1.10						Prompting	0.376	0.478	18.70	4.73	1.66
ACON <u>UT</u> 74.74 13.13 4.93 3.85						ACON <u>UT</u>	0.373	0.494	17.14	4.71	1.57
ACON <u>UTCO</u> 72.63 11.54 4.54 1.91						ACON UTCO	0.335	0.458	17.79	4.65	1.50
Observation Con	pression					Observation Con	npression	1			
LLMLingua	71.58	11.89	7.38	6.14		LLMLingua	0.320	0.414	14.23	5.16	1.35
Prompting	55.79	12.24	6.44	2.68		Prompting	0.288	0.397	11.64	3.41	0.45
ACON UT	73.68	10.83	6.55	3.85		ACON UT	0.364	0.475	16.33	4.97	1.28
ACON UTCO	72.63	10.28	6.17	2.88		ACON UTCO	0.336	0.461	14.00	4.22	0.81
	AppV	Vorld			О	fficeBench		8	3-objecti	ve QA	
	_	Te	acher	80		77.89	0.3	9	0.386		
≥ 50 + 50.0	0		action			75.79	1 0.5	٦	0.366		
96 96	47.00			75 - Tea	cher_	- 73.68	0.3	ο		0.3	79
. I		44.60		/3 7		- /3.08 -					
45 - 05 do			2.30	70	.53		Teach 0.3		(.370	-
,		4	2.30	70			[0.3				
	Sı	mall Mod	lels: 🔲	Qwen3-1	4B	Qwen3-8B	Phi-4	gpt-4.1	-mini		

Figure 4: **Results of distilled compressors on history compression** with gpt-4.1 as the agent. Grey dotted lines denote performance using the gpt-4.1 teacher compressor. Student models (Qwen3-14B, Qwen3-8B, Phi-4) are distilled from gpt-4.1 compressor using the optimized compression guideline after UT step, and evaluated across all benchmarks. We also include gpt-4.1-mini without distillation, showing that even a small model can serve as an effective compressor without additional training.

Applying only the utility maximization step (UT) improves performance while reducing token cost across all benchmarks, whereas the compression maximization step (CO) further lowers token cost but may slightly hurt accuracy, except in AppWorld where it even yields additional gains.

4.3 Compressor Distillation

We distill the compressor with optimized guidelines after UT step into smaller models such as Qwen3-14B, Qwen3-8B (Yang et al., 2025), and Phi-4 (Abdin et al., 2024) using LoRA (Hu et al., 2022). As shown in Figure 4, the distilled compressors retain over 95% of the performance of the gpt-4.1 teacher (indicated by the grey dotted line) while reducing computational overhead. We also observe that gpt-4.1-mini, even without any distillation, can serve as an effective lightweight compressor on OfficeBench and QA. This indicates that small models can reliably replace large LLM-based compressors when equipped with optimized guidelines. These results confirm that small models are sufficient for compression, enabling the expensive LLM to be only reserved for the agent.

4.4 ACON FOR DISTILLED SMALL AGENTS

We examine whether ACON also benefits smaller LLM agents, which are particularly vulnerable to long-horizon inefficiency. Without compression, models such as Qwen3-14B often fail on medium and hard tasks due to distracting context. As shown in Figure 5, ACON substantially improves their performance: on AppWorld, Qwen3-14B improves from 26.8% to 33.9%, and on 8-objective QA from 0.158 to 0.197 EM. These results demonstrate that ACON acts as an equalizer, enabling smaller agents with concise but informative contexts to approach the performance of larger models.

4.5 Analysis

Compression threshold: moderate value yields the best trade-off. In Figure 6, we provide ablations on threshold for compression in Equation 3 and Equation 4. Results show that smaller

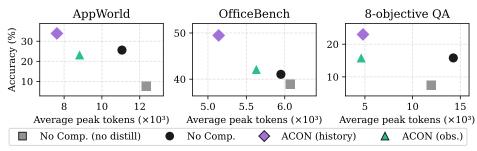


Figure 5: **Performance-efficiency trade-off of the Qwen3-14B agent** distilled from gpt-4.1 trajectories. For distilled compressors, we use the same distillation setting as in Figure 4. Compared to the baseline without compression, our framework ACON provides compressed trajectories combined with a distilled compressor, enabling the distilled agent to achieve consistently higher accuracy while requiring substantially fewer peak input tokens across all benchmarks.

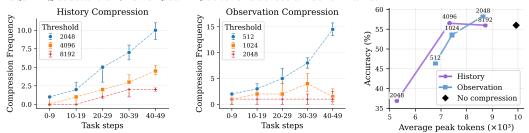


Figure 6: **Ablation studies on thresholds for compression** on AppWorld with gpt-4.1. (1) the number of compressions (compression frequency) for each length of task trajectories (task steps). (2) the performance comparison for each threshold setting.

thresholds reduce tokens but incur more frequent compression calls and degrade accuracy, while larger thresholds preserve accuracy with higher cost. Moderate values (4096 for history, 1024 for observation) provide the best trade-off, maintaining accuracy close to no compression while still reducing peak tokens substantially.

Optimization cost. It is important to clarify the cost of the guideline optimization introduced in Section 3.2, as it may initially appear expensive to practitioners. We compute all costs using the official API pricing described in Section B.2, and the results are summarized in Table 4. As shown, **the total cost is around \$20 for an AppWorld-scale domain.** The optimization consists of two phases: (1) collecting trajectories with the agent model and (2) optimizing the guideline with the optimizer model. For the rollout phase, we use fewer than 100 training examples for all benchmarks used (Appendix B.1).

Table 4: API cost for compression guideline optimization on AppWorld benchmark with 90 training examples.

Stage (LM)	Cost (\$)
Rollout	
without comp. (qpt-4.1)	10.7
with comp. (gpt-4.1)	8.9
Optimization	
UT step (03)	1.2
\overline{CO} step $(\circ 3)$	0.7
Total	21.5

Executing both compressed and uncompressed trajectories

with gpt-4.1 costs around \$20 for 90 examples. This cost can be amortized by integrating rollout collection into normal online task execution (Wang et al., 2025c), and it can be reduced further by using fewer examples when needed. Notably, we require at most two trajectories per example, which is substantially lower than reinforcement learning based approaches such as GRPO that require multiple trajectories per example for advantage estimation (Shao et al., 2024). For the optimization phase, we generate five candidate prompts per iteration and perform one \overline{UT} step and one \overline{CO} step. The total cost of this stage is under \$2 with the 03 model. Since our method does not require any model weight updates or GPU resources, it remains practical and easily applicable to both open and proprietary models via API access.

Prompt optimizer: o3 + contrastive feedback works best. We analyze how the choice of optimizer and the use of contrastive feedback affect compression guideline quality. As shown in Table 3, the default o3 with contrastive feedback yields the best performance, while removing contrastive feedback (only using failed trajectories) or switching to other models results in lower accuracy. Although o3 shows the best performance, we also demonstrate that the optimizer model can be re-

Table 3: **Ablation studies on the prompt optimizer** in AppWorld, gpt-4.1 agent and history compressor. Default is o3 optimizer with task contrastive feedback.

Optimizer model	Task contrastive	Average Acc.
о3	\checkmark	51.2
03	Х	50.6 (-0.6)
gpt-4.1	\checkmark	47.6 (-3.6)
gpt-5	\checkmark	50.6 (-0.6)

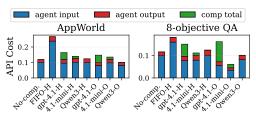


Figure 7: **API cost comparison.** H/O denote history/observation compression. Listed models are compressors; the agent is gpt-4.1.

placed to weaker models such as gpt-4.1, showing it still yields sufficiently fine guideline compared to the baseline guideline.

Cost analysis with an API cost proxy. An important point of discussion is the extent to which ACON improves efficiency. To clarify, ACON primarily targets *memory efficiency* (*e.g.*, reducing peak tokens and dependency), rather than providing overall API or computational cost reductions. Our experiments show that ACON reliably reduces memory cost while preserving or even improving task performance. Nonetheless, to prevent misunderstandings and to better inform future work, we also analyze whether ACON reduces end-to-end computational cost. We use API cost as a proxy for computational overhead, accounting for discounted cached tokens as described in Section B.2.

We report the API cost analysis in Figure 7, using gpt-4.1 as the agent and gpt-4.1, gpt-4.1-mini, and Qwen3-14B as compressors. The analysis yields three observations. First, using small models reduces compressor cost: replacing gpt-4.1 with smaller models such as gpt-4.1-mini or Qwen3-14B lowers overhead while preserving compression quality. Second, observation compression reduces total API cost: in AppWorld, using Qwen3-14B compressor lowers the per-

Table 5: Effect of compression threshold on peak tokens and API cost.

Thr.	Peak (10 ³)	Agent (\$)	Comp. (\$)
2048	5.36	.1440	.0720
4096	7.33	.1253	.0360
8192	8.70	.1179	.0140

task cost by 15% (from \$0.119 to \$0.101) because observations are compressed before being cached in the agent. Third, *history compression may increase or decrease total cost depending on the task*. In AppWorld, cost increases by 4% (from \$0.119 to \$0.124) with Qwen3-14B compressor due to KV-cache invalidation and occasional increases in task length. In contrast, in 8-objective QA, cost decreases by 4% (from \$0.116 to \$0.112) with gpt-4.1-mini compressor, as fewer compression events lead to minimal KV-cache disruption while still shortening inputs. The benefit of history compression therefore depends on how often the compressed history is reused, which is controlled by the compression threshold. As shown in Figure 6 and Table 5, the threshold determines compression frequency and trades off peak tokens against total API cost. We recommend tuning it on a validation set to achieve the desired balance.

We include a more detailed discussion of the limitation, together with practical guidance on threshold selection and compressor choice, in Appendix A.

5 CONCLUSION

We presented **Agent Context Optimization** (**ACON**), a unified framework that systematically compresses both interaction histories and environment observations for long-horizon LLM agents. Unlike prior work that relies on naive prompting or narrow domains, ACON introduces compression guideline optimization in natural language space, enabling adaptive and model-agnostic compression. Experiments on AppWorld, OfficeBench, and Multi-objective QA show that ACON reduces peak tokens by 26-54% while maintaining or even improving task success. Beyond memory efficiency, we demonstrate that optimized compressors can be distilled into smaller models, substantially lowering overhead without sacrificing performance. Moreover, by supplying concise yet informative contexts, ACON allows small agents such as Qwen3-14B to approach the performance of much larger models. Overall, our findings highlight that ACON lays a foundation for more general, cost-effective, and deployable long-horizon LLM agents.

ETHICS STATEMENT

This work investigates optimized context compression framework for long-horizon LLM agents. It **does not** involve human subjects, user studies, or the collection of personally identifiable information. All experiments are conducted on **publicly available** benchmarks released under their respective licenses, which, to the best of our knowledge, do not contain sensitive personal data.

REPRODUCIBLE STATEMENTS

We conduct all experiments using the Azure OpenAI endpoint with fixed model snapshots and seed, specifically gpt-4.1-2025-04-14 and gpt-4.1-mini-2025-04-14, to ensure reproducibility. Our implementation relies on PyTorch (Paszke et al., 2019), the Hugging Face Transformers library (Wolf et al., 2020), and the TRL library for training and vLLM (Kwon et al., 2023) for inference. Additional implementation details are provided in the Appendix B. We will also release our codebase to enable the research community to fully reproduce, verify, and extend our work on long-horizon LLM agents.

REFERENCES

Marah I Abdin, Jyoti Aneja, Harkirat S. Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J. Hewett, Mojan Javaheripi, Piero Kauffmann, James R. Lee, Yin Tat Lee, Yuanzhi Li, Weishung Liu, Caio C. T. Mendes, Anh Nguyen, Eric Price, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital Shah, Xin Wang, Rachel Ward, Yue Wu, Dingli Yu, Cyril Zhang, and Yi Zhang. Phi-4 technical report. *arXiv*, 2412.08905, 2024. URL https://doi.org/10.48550/arXiv.2412.08905.

Anthropic. Introducing the model context protocol. https://www.anthropic.com/news/model-context-protocol, 2024a.

Anthropic. The claude 3 model family: Opus, sonnet, haiku. 2024b. URL https://api.semanticscholar.org/CorpusID:268232499.

Rogerio Bonatti, Dan Zhao, Francesco Bonacci, Dillon Dupont, Sara Abdali, Yinheng Li, Yadong Lu, Justin Wagle, Kazuhito Koishida, Arthur Bucker, Lawrence Jang, and Zack Hui. Windows agent arena: Evaluating multi-modal OS agents at scale. *arXiv*, 2409.08264, 2024. URL https://doi.org/10.48550/arXiv.2409.08264.

Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit S. Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke Marris, Sam Petulla, Colin Gaffney, Asaf Aharoni, Nathan Lintz, Tiago Cardal Pais, Henrik Jacobsson, Idan Szpektor, Nan-Jiang Jiang, Krishna Haridasan, Ahmed Omran, Nikunj Saunshi, Dara Bahri, Gaurav Mishra, Eric Chu, Toby Boyd, Brad Hekman, Aaron Parisi, Chaoyi Zhang, Kornraphop Kawintiranon, Tania Bedrax-Weiss, Oliver Wang, Ya Xu, Ollie Purkiss, Uri Mendlovic, Ilaï Deutel, Nam Nguyen, Adam Langley, Flip Korn, Lucia Rossazza, Alexandre Ramé, Sagar Waghmare, Helen Miller, Nathan Byrd, Ashrith Sheshan, Raia Hadsell Sangnie Bhardwaj, Pawel Janus, Tero Rissa, Dan Horgan, Sharon Silver, Ayzaan Wahid, Sergey Brin, Yves Raimond, Klemen Kloboves, Cindy Wang, Nitesh Bharadwaj Gundavarapu, Ilia Shumailov, Bo Wang, Mantas Pajarskas, Joe Heyward, Martin Nikoltchev, Maciej Kula, Hao Zhou, Zachary Garrett, Sushant Kafle, Sercan Arik, Ankita Goel, Mingyao Yang, Jiho Park, Koji Kojima, Parsa Mahmoudieh, Koray Kavukcuoglu, Grace Chen, Doug Fritz, Anton Bulyenov, Sudeshna Roy, Dimitris Paparas, Hadar Shemtov, Bo-Juen Chen, Robin Strudel, David Reitter, Aurko Roy, Andrey Vlasov, Changwan Ryu, Chas Leichner, Haichuan Yang, Zelda Mariet, Denis Vnukov, Tim Sohn, Amy Stuart, Wei Liang, Minmin Chen, Praynaa Rawlani, Christy Koh, JD Co-Reyes, Guangda Lai, Praseem Banzal, Dimitrios Vytiniotis, Jieru Mei, and Mu Cai. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. arXiv, 2507.06261, 2025. URL https://doi.org/10.48550/arXiv.2507.06261.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu,

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

Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samual Stevens, Boshi Wang, Huan Sun, and Yu Su. Mind2web: Towards a generalist agent for the web. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/5950bf290a1570ea401bf98882128160-Abstract-Datasets_and_Benchmarks.html.

Dongge Han, Menglin Xia, Daniel Madrigal Diaz, Samuel Kessler, Ankur Mallick, Xuchao Zhang, Mirian del Carmen Hipolito Garcia, Jin Xu, Victor Rühle, and Saravan Rajmohan. Enhancing reasoning capabilities of small language models with blueprints and prompt template search. *arXiv*, 2506.08669, 2025. URL https://doi.org/10.48550/arXiv.2506.08669.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=nZeVKeeFYf9.

Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. Llmlingua: Compressing prompts for accelerated inference of large language models. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pp. 13358–13376. Association for Computational Linguistics, 2023. URL https://doi.org/10.18653/v1/2023.emnlp-main.825.

Huiqiang Jiang, Qianhui Wu, Xufang Luo, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. Longllmlingua: Accelerating and enhancing llms in long context scenarios via prompt compression. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, *ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pp. 1658–1677. Association for Computational Linguistics, 2024. URL https://doi.org/10.18653/v1/2024.acl-long.91.

- 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, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net, 2024. URL https://openreview.net/forum?id=VTF8yNQM66.
- Bowen Jin, Hansi Zeng, Zhenrui Yue, Dong Wang, Hamed Zamani, and Jiawei Han. Search-r1: Training llms to reason and leverage search engines with reinforcement learning. *arXiv*, 2503.09516, 2025. URL https://doi.org/10.48550/arXiv.2503.09516.
- Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan A, Saiful Haq, Ashutosh Sharma, Thomas T. Joshi, Hanna Moazam, Heather Miller, Matei Zaharia, and Christopher Potts. DSPy: Compiling declarative language model calls into state-of-the-art pipelines. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=sy5N0zy5Od.
- Minsoo Kim, Arnav Kundu, Han-Byul Kim, Richa Dixit, and Minsik Cho. Epicache: Episodic kv cache management for long conversational question answering, 2025. URL https://arxiv.org/abs/2509.17396.
- Yoon Kim and Alexander M. Rush. Sequence-level knowledge distillation. In Jian Su, Xavier Carreras, and Kevin Duh (eds.), *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, pp. 1317–1327. The Association for Computational Linguistics, 2016. URL https://doi.org/10.18653/v1/d16-1139.
- Thomas Kwa, Ben West, Joel Becker, Amy Deng, Katharyn Garcia, Max Hasin, Sami Jawhar, Megan Kinniment, Nate Rush, Sydney von Arx, Ryan Bloom, Thomas Broadley, Haoxing Du, Brian Goodrich, Nikola Jurkovic, Luke Harold Miles, Seraphina Nix, Tao Lin, Neev Parikh, David Rein, Lucas Jun Koba Sato, Hjalmar Wijk, Daniel M. Ziegler, Elizabeth Barnes, and Lawrence Chan. Measuring AI ability to complete long tasks. *arXiv*, 2503.14499, 2025. URL https://doi.org/10.48550/arXiv.2503.14499.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur P. Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. Natural questions: a benchmark for question answering research. *Trans. Assoc. Comput. Linguistics*, 7:452–466, 2019. URL https://doi.org/10.1162/tacl_a_00276.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In Jason Flinn, Margo I. Seltzer, Peter Druschel, Antoine Kaufmann, and Jonathan Mace (eds.), *Proceedings of the 29th Symposium on Operating Systems Principles, SOSP 2023, Koblenz, Germany, October 23-26, 2023*, pp. 611–626. ACM, 2023. URL https://doi.org/10.1145/3600006.3613165.
- Dongjun Lee, Juyong Lee, Kyuyoung Kim, Jihoon Tack, Jinwoo Shin, Yee Whye Teh, and Kimin Lee. Learning to contextualize web pages for enhanced decision making by LLM agents. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025.* OpenReview.net, 2025. URL https://openreview.net/forum?id=3Gzz7ZQLiz.
- Yucheng Li, Bo Dong, Frank Guerin, and Chenghua Lin. Compressing context to enhance inference efficiency of large language models. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pp. 6342–6353. Association for Computational Linguistics, 2023. URL https://doi.org/10.18653/v1/2023.emnlp-main.391.
- Adyasha Maharana, Dong-Ho Lee, Sergey Tulyakov, Mohit Bansal, Francesco Barbieri, and Yuwei Fang. Evaluating very long-term conversational memory of llm agents. *arXiv preprint arXiv:2402.17753*, 2024.

```
OpenAI. Introducing chatgpt agent: bridging research and action. https://openai.com/index/introducing-chatgpt-agent/, 2025a.
```

- OpenAI. Introducing gpt-4.1 in the api. https://openai.com/index/gpt-4-1/, 2025b.
- OpenAI. Introducing gpt-5. https://openai.com/index/introducing-gpt-5/, 2025c.
- OpenAI. Introducing openai o3 and o4-mini. https://openai.com/index/introducing-o3-and-o4-mini/, 2025d.
- Charles Packer, Vivian Fang, Shishir_G Patil, Kevin Lin, Sarah Wooders, and Joseph_E Gonzalez. Memgpt: Towards llms as operating systems. 2023.
- Zhuoshi Pan, Qianhui Wu, Huiqiang Jiang, Menglin Xia, Xufang Luo, Jue Zhang, Qingwei Lin, Victor Rühle, Yuqing Yang, Chin-Yew Lin, H. Vicky Zhao, Lili Qiu, and Dongmei Zhang. Llmlingua-2: Data distillation for efficient and faithful task-agnostic prompt compression. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics*, *ACL 2024*, *Bangkok*, *Thailand and virtual meeting*, *August 11-16*, 2024, pp. 963–981. Association for Computational Linguistics, 2024. URL https://doi.org/10.18653/v1/2024.findings-acl.57.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Neural Information Processing Systems (NeurIPS)*, 2019.
- Reid Pryzant, Dan Iter, Jerry Li, Yin Tat Lee, Chenguang Zhu, and Michael Zeng. Automatic prompt optimization with "gradient descent" and beam search. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pp. 7957–7968. Association for Computational Linguistics, 2023. URL https://doi.org/10.18653/v1/2023.emnlp-main.494.
- Minju Seo, Jinheon Baek, Seongyun Lee, and Sung Ju Hwang. Efficient long context language model retrieval with compression. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2025, Vienna, Austria, July 27 August 1, 2025*, pp. 15251–15268. Association for Computational Linguistics, 2025. URL https://aclanthology.org/2025.acl-long.740/.
- Shivam Shandilya, Menglin Xia, Supriyo Ghosh, Huiqiang Jiang, Jue Zhang, Qianhui Wu, Victor Rühle, and Saravan Rajmohan. TACO-RL: task aware prompt compression optimization with reinforcement learning. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), Findings of the Association for Computational Linguistics, ACL 2025, Vienna, Austria, July 27 August 1, 2025, pp. 1582–1597. Association for Computational Linguistics, 2025. URL https://aclanthology.org/2025.findings-acl.81/.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv*, 2402.03300, 2024. URL https://doi.org/10.48550/arXiv.2402.03300.
- Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H. Chi, Nathanael Schärli, and Denny Zhou. Large language models can be easily distracted by irrelevant context. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 31210–31227. PMLR, 2023. URL https://proceedings.mlr.press/v202/shi23a.html.
- Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew J. Hausknecht. Alfworld: Aligning text and embodied environments for interactive learning. In

758

759

760

761

762 763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788 789

790

791

792

793

794

796

797

798

799 800

801

802 803

804

805

806

807

808

9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021. URL https://openreview.net/forum?id= 0IOX0YcCdTn.

Calvin Smith. Openhands context condensensation for more efficient ai agents. *All Hands AI Blog*, April 2025. URL https://www.all-hands.dev/blog/openhands-context-condensensation-for-more-efficient-ai-agents.

Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, second edition, 2018.

Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, Zhuofu Chen, Jialei Cui, Hao Ding, Mengnan Dong, Angang Du, Chenzhuang Du, Dikang Du, Yulun Du, Yu Fan, Yichen Feng, Kelin Fu, Bofei Gao, Hongcheng Gao, Peizhong Gao, Tong Gao, Xinran Gu, Longyu Guan, Haiging Guo, Jianhang Guo, Hao Hu, Xiaoru Hao, Tianhong He, Weiran He, Wenyang He, Chao Hong, Yangyang Hu, Zhenxing Hu, Weixiao Huang, Zhiqi Huang, Zihao Huang, Tao Jiang, Zhejun Jiang, Xinyi Jin, Yongsheng Kang, Guokun Lai, Cheng Li, Fang Li, Haoyang Li, Ming Li, Wentao Li, Yanhao Li, Yiwei Li, Zhaowei Li, Zheming Li, Hongzhan Lin, Xiaohan Lin, Zongyu Lin, Chengyin Liu, Chenyu Liu, Hongzhang Liu, Jingyuan Liu, Junqi Liu, Liang Liu, Shaowei Liu, T. Y. Liu, Tianwei Liu, Weizhou Liu, Yangyang Liu, Yibo Liu, Yiping Liu, Yue Liu, Zhengying Liu, Enzhe Lu, Lijun Lu, Shengling Ma, Xinyu Ma, Yingwei Ma, Shaoguang Mao, Jie Mei, Xin Men, Yibo Miao, Siyuan Pan, Yebo Peng, Ruoyu Qin, Bowen Qu, Zeyu Shang, Lidong Shi, Shengyuan Shi, Feifan Song, Jianlin Su, Zhengyuan Su, Xinjie Sun, Flood Sung, Heyi Tang, Jiawen Tao, Qifeng Teng, Chensi Wang, Dinglu Wang, Feng Wang, Haiming Wang, Jianzhou Wang, Jiaxing Wang, Jinhong Wang, Shengjie Wang, Shuyi Wang, Yao Wang, Yejie Wang, Yiqin Wang, Yuxin Wang, Yuzhi Wang, Zhaoji Wang, Zhengtao Wang, Zhexu Wang, Chu Wei, Qianqian Wei, Wenhao Wu, Xingzhe Wu, Yuxin Wu, Chenjun Xiao, Xiaotong Xie, Weimin Xiong, Boyu Xu, Jing Xu, Jinjing Xu, L. H. Xu, Lin Xu, Suting Xu, Weixin Xu, Xinran Xu, Yangchuan Xu, Ziyao Xu, Junjie Yan, Yuzi Yan, Xiaofei Yang, Ying Yang, Zhen Yang, Zhilin Yang, Zonghan Yang, Haotian Yao, Xingcheng Yao, Wenjie Ye, Zhuorui Ye, Bohong Yin, Longhui Yu, Enming Yuan, Hongbang Yuan, Mengjie Yuan, Haobing Zhan, Dehao Zhang, Hao Zhang, Wanlu Zhang, Xiaobin Zhang, Yangkun Zhang, Yizhi Zhang, Yongting Zhang, Yu Zhang, Yutao Zhang, Yutong Zhang, Zheng Zhang, Haotian Zhao, Yikai Zhao, Huabin Zheng, Shaojie Zheng, Jianren Zhou, Xinyu Zhou, Zaida Zhou, Zhen Zhu, Weiyu Zhuang, and Xinxing Zu. Kimi k2: Open agentic intelligence, 2025. URL https://arxiv.org/abs/2507.20534.

Harsh Trivedi, Tushar Khot, Mareike Hartmann, Ruskin Manku, Vinty Dong, Edward Li, Shashank Gupta, Ashish Sabharwal, and Niranjan Balasubramanian. Appworld: A controllable world of apps and people for benchmarking interactive coding agents. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pp. 16022–16076. Association for Computational Linguistics, 2024. URL https://doi.org/10.18653/v1/2024.acl-long.850.

Qingyue Wang, Yanhe Fu, Yanan Cao, Shuai Wang, Zhiliang Tian, and Liang Ding. Recursively summarizing enables long-term dialogue memory in large language models. *Neurocomputing*, 639:130193, 2025a.

Weixuan Wang, Dongge Han, Daniel Madrigal Diaz, Jin Xu, Victor Rühle, and Saravan Rajmohan. Odysseybench: Evaluating Ilm agents on long-horizon complex office application workflows, 2025b. URL https://arxiv.org/abs/2508.09124.

Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji. Executable code actions elicit better LLM agents. In *Forty-first International Conference on Machine Learning, ICML* 2024, *Vienna, Austria, July* 21-27, 2024. OpenReview.net, 2024a. URL https://openreview.net/forum?id=jJ9BoXAfFa.

Zilong Wang, Yuedong Cui, Li Zhong, Zimin Zhang, Da Yin, Bill Yuchen Lin, and Jingbo Shang. Officebench: Benchmarking language agents across multiple applications for office automation. *arXiv*, 2407.19056, 2024b. URL https://doi.org/10.48550/arXiv.2407.19056.

- Zora Zhiruo Wang, Jiayuan Mao, Daniel Fried, and Graham Neubig. Agent workflow memory. In *Forty-second International Conference on Machine Learning*, 2025c. URL https://openreview.net/forum?id=NTAhi2JEEE.
- Jason Wei, Zhiqing Sun, Spencer Papay, Scott McKinney, Jeffrey Han, Isa Fulford, Hyung Won Chung, Alex Tachard Passos, William Fedus, and Amelia Glaese. Browsecomp: A simple yet challenging benchmark for browsing agents. *arXiv*, 2504.12516, 2025. URL https://doi.org/10.48550/arXiv.2504.12516.
- Jeffrey Willette, Heejun Lee, Youngwan Lee, Myeongjae Jeon, and Sung Ju Hwang. Training free exponential context extension via cascading KV cache. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025.* OpenReview.net, 2025. URL https://openreview.net/forum?id=dSneEp59yX.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In *EMNLP 2020 Demos, Online, November 16-20, 2020*, pp. 38–45, 2020. URL https://doi.org/10.18653/v1/2020.emnlp-demos.6.
- Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming language models with attention sinks. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net, 2024. URL https://openreview.net/forum?id=NG7sS51zVF.
- Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu, Yiheng Xu, Shuyan Zhou, Silvio Savarese, Caiming Xiong, Victor Zhong, and Tao Yu. Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 15, 2024, 2024. URL http://papers.nips.cc/paper_files/paper/2024/hash/5d413e48f84dc61244b6be550f1cd8f5-Abstract-Datasets_and_Benchmarks_Track.html.
- Fangyuan Xu, Weijia Shi, and Eunsol Choi. RECOMP: improving retrieval-augmented lms with context compression and selective augmentation. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL https://openreview.net/forum?id=mlJLVigNHp.
- Wujiang Xu, Zujie Liang, Kai Mei, Hang Gao, Juntao Tan, and Yongfeng Zhang. A-MEM: agentic memory for LLM agents. *arXiv*, 2502.12110, 2025. URL https://doi.org/10.48550/arXiv.2502.12110.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jian Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report. arXiv, 2505.09388, 2025. URL https://doi.org/10.48550/arXiv.2505.09388.
- Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V. Le, Denny Zhou, and Xinyun Chen. Large language models as optimizers. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024a. URL https://openreview.net/forum?id=Bb4VGOWELI.

- 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. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 15, 2024, 2024b. URL http://papers.nips.cc/paper_files/paper/2024/hash/5a7c947568c1b1328ccc5230172e1e7c-Abstract-Conference.html.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik R. Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL https://openreview.net/forum?id=WE_vluYUL-X.
- Chanwoong Yoon, Taewhoo Lee, Hyeon Hwang, Minbyul Jeong, and Jaewoo Kang. Compact: Compressing retrieved documents actively for question answering. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pp. 21424–21439. Association for Computational Linguistics, 2024. URL https://doi.org/10.18653/v1/2024.emnlp-main.1194.
- Mert Yüksekgönül, Federico Bianchi, Joseph Boen, Sheng Liu, Pan Lu, Zhi Huang, Carlos Guestrin, and James Zou. Optimizing generative AI by backpropagating language model feedback. *Nature*, 639(8055):609–616, 2025. URL https://doi.org/10.1038/s41586-025-08661-4.
- Jintian Zhang, Yuqi Zhu, Mengshu Sun, Yujie Luo, Shuofei Qiao, Lun Du, Da Zheng, Huajun Chen, and Ningyu Zhang. Lightthinker: Thinking step-by-step compression. *arXiv*, 2502.15589, 2025. URL https://doi.org/10.48550/arXiv.2502.15589.
- Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. Webarena: A realistic web environment for building autonomous agents. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL https://openreview.net/forum?id=oKn9c6ytLx.
- Zijian Zhou, Ao Qu, Zhaoxuan Wu, Sunghwan Kim, Alok Prakash, Daniela Rus, Jinhua Zhao, Bryan Kian Hsiang Low, and Paul Pu Liang. MEM1: learning to synergize memory and reasoning for efficient long-horizon agents. *arXiv*, 2506.15841, 2025. URL https://doi.org/10.48550/arXiv.2506.15841.

A LIMITATIONS & FUTURE WORKS

Our work addresses the context management problem in long-horizon LLM agents and proposes a framework for optimized context compression. While the method effectively reduces token costs with minimal performance degradation, it also presents several limitations.

A primary limitation is computational overhead. As discussed in Section 4.5, history compression can in some cases slightly increase total cost, since additional steps may be required to complete challenging tasks with the compressed history. Moreover, it breaks the KV-cache of transformer-based LLMs, which forces re-computation of compressed histories. This effect is particularly evident in the FIFO baseline. Observation compression alleviates some of this overhead, but the generative compression procedure itself introduces latency, slowing down agent response time (Lee et al., 2025). A promising future direction is the development of KV-cache–level compression or eviction strategies. Prior work has considered KV-cache compression for single-step reasoning (Zhang et al., 2025), long documents (Xiao et al., 2024; Willette et al., 2025), and long conversations (Kim et al., 2025). However, its role in multi-turn, long-horizon agents remains underexplored.

Another limitation is model and benchmark coverage. Our experiments primarily evaluate GPT models due to budgetary constraints. Although the framework is designed to be model-agnostic, its generalizability to other foundation models such as Gemini (Comanici et al., 2025) or Claude (Anthropic, 2024b) remains unverified. Similarly, we were unable to include large-scale open-source models such as DeepSeek-R1 (DeepSeek-AI et al., 2025) or Qwen3-235B (Yang et al., 2025) due to limited resources. Extending the analysis to these models would provide stronger evidence of robustness and broaden the applicability of our conclusions. We also validated our approach on three benchmarks. While these benchmarks reflect realistic agent settings, full real-world deployability remains an open question. Evaluating the method in real-world environments, where tasks are less controlled and constraints are more complex, would be a valuable direction for future work.

Furthermore, the optimization of the compression guideline also has limitations. Our method, along with prior prompt optimization approaches (Yüksekgönül et al., 2025; Khattab et al., 2024; Pryzant et al., 2023; Yang et al., 2024a), supports the view that updating natural language instructions using LLM-generated feedback is a valid strategy for improving LLM systems. However, unlike numerical gradient-based methods, this optimization provides no convergence guarantee. We partially address this issue by sampling multiple candidate guidelines and selecting one during training, but this remains a heuristic solution. A deeper analysis of the optimization process and more principled methods for optimizing the objective in Section 3.2 would be valuable directions for future work.

Another limitation concerns distillation quality. Although our distilled models retain most of the teacher's behavior, they do not achieve perfectly identical performance. We expect that this gap can be further reduced by increasing the amount of training data beyond the 100 examples per domain used in this work.

B EXPERIMENTAL SETUP DETAILS

B.1 DATASETS

AppWorld (Trivedi et al., 2024). AppWorld is our primary benchmark, providing a high-quality execution environment that integrates nine everyday applications (Spotify, SimpleNote, Amazon, Venmo, Gmail, Splitwise, File system, Todoist, and Phone) through 457 APIs. It also includes realistic simulations of approximately 100 functional users. This benchmark is particularly suitable for evaluating long-horizon productivity agents, as its multi-step tasks require an average of 42.5 API calls per task. We follow the official split, using 90 training tasks for guideline optimization and distillation, and 168 test-normal tasks for evaluation. An example trajectory from AppWorld is provided in Example E.1.

OfficeBench (Wang et al., 2024b). OfficeBench is a benchmark for office automation using applications such as Word, Excel, PDF, Calendar, Email, Shell, and Calculator. It evaluates the ability of agents to coordinate across multiple apps to complete complex tasks, making it well suited for long-horizon scenarios. Tasks are categorized as 1-app, 2-app, or 3-app depending on the number of applications required. We restrict our experiments to text-related tasks, excluding those requiring

Table 6: Example tasks across benchmarks.

Benchmark / Difficulty	Example Task
AppWorld	
Easy	Mark "Taking a solo backpacking trip" in my Bucket List Simple Note as not done.
Medium	Like all the Venmo transactions from today involving any of my roommates on my Venmo social feed.
Hard	Start playing a playlist on Spotify that has enough songs for my work- out today. I do not want to have to change the playlist in the middle of my workout. The workout plan is in Simple Note.
OfficeBench	
1-app	Create a new Word file called random_paragraph.docx and add the content in random_paragraph.txt to it.
2-app	Analyze Excel data of students' grade and generate a teaching report in teaching.docx.
3-app	Read company revenues and send an email with subject revenues, containing data to Bob for reporting, also write a revenue.docx to summarize it.
8-objective QA	
_	who wrote the song <i>Oceans Where Feet May Fail</i> ?; who plays Eddie the Eagle in the movie?; when was the last time England were in the final of World Cup?; who plays Chelsea's mom on <i>Young and the Restless</i> ?; what is the largest coin in the US?; who sang <i>Even the Bad Times Are Good</i> ?; who sings <i>This Is My Town</i> country song?; which of the Guianas is not an independent country?

OCR, as OCR quality could confound the evaluation. Since no official split is available, we randomly partition the tasks into training and test sets with a 1:1 ratio, resulting in 92 training tasks and 95 test tasks. We additionally refine the dataset by removing ambiguous tasks and ensuring that synthetic files (testbeds) are not shared across splits.

8-Objective QA (Zhou et al., 2025). The 8-objective QA benchmark simulates deep researchstyle agentic tasks. Unlike conventional multi-hop QA, which requires answering a single question using multiple pieces of evidence, this benchmark poses eight distinct questions within one task, and the agent must provide answers to all of them at the end. This design creates a more challenging setting for long-horizon agents. Following Zhou et al. (2025), we construct each task by grouping eight questions together. Questions are drawn from NaturalQuestions (Kwiatkowski et al., 2019), resulting in 100 training tasks (from the train split) and 100 test tasks (from the test split). For retrieval, we use a BM25 retriever over the 2018 Wikipedia knowledge base, following Jin et al. (2025).

We include the example task of each benchmark in Table 6.

B.2 EVALUATION METRICS

For efficiency evaluation, we adopt two metrics—peak tokens and dependency—introduced in Light-Thinker (Zhang et al., 2025) and MEM1 (Zhou et al., 2025).

Peak tokens. Peak tokens are measured as the maximum number of tokens observed in any single sequence throughout the agent's trajectory, excluding system prompts. This metric serves as a proxy for inference-time memory requirements and corresponds to the maximum peak shown in Figure 2.

Dependency. Dependency is defined as the area under the curve in Figure 2. At each step t, given the number of input tokens $n_i^{(t)}$ and output tokens $n_o^{(t)}$, it is calculated as:

$$\sum_{t \in [T]} \frac{(n_i^{(t)} + 2n_o^{(t)}) \times n_o^{(t)}}{2}.$$
(11)

This metric approximates the cumulative computational cost incurred by action generation across the trajectory.

API Cost. For the cost analysis in Section 4.5, we use the official OpenAI pricing (as of September 2025) for gpt-4.1 and gpt-4.1-mini (OpenAI, 2025b). Specifically, gpt-4.1 is priced at \$3.00 per 1M input tokens, \$0.75 per 1M cached input tokens, and \$12.00 per 1M output tokens. For gpt-4.1-mini, the costs are \$0.80 per 1M input tokens, \$0.20 per 1M cached input tokens, and \$3.20 per 1M output tokens. For Qwen3-14B (Yang et al., 2025), since no official API pricing is available, we approximate the cost using OpenRouter¹: \$0.06 per 1M input tokens, \$0.015 per 1M cached input tokens, and \$0.24 per 1M output tokens.

B.3 IMPLEMENTATION DETAILS & HYPERPARAMETERS

API Inference. We set temperature 0.0 and fix the seed 42. Note that there is still non-determinism with fixing the seed and setting temperature as 0. To reduce the instability, we use the API snapshot form Azure OpenAI endpoint gpt-4.1-2025-04-14 and gpt-4.1-mini-2025-04-14.

Compression. For history compression, we set $T_{\rm hist}=4096$ for AppWorld and OfficeBench, and 2048 for 8-objective QA. We keep the last action, observation pair to preserve the latest information. This is the same for ACON and all baselines. For observation compression, we set $T_{\rm obs}=1024$ for AppWorld, 512 for OfficeBench, and 400 for 8-objective QA.

Prompt Optimization. We use the OpenAI o3 model (OpenAI, 2025d) for both analysis and update of prompts. During the update stage, we sample 5 candidate prompts and select the one that performs best on a subset of the training set.

The prompts used in each stage and step are provided as follows:

- Analysis prompt for <u>UT</u> step: Prompt E.1
- Update prompt for UT step: Prompt E.2
- Analysis prompt for \overline{CO} step: Prompt E.3
- Update prompt for \overline{CO} step: Prompt E.4

We also provide the detailed procedure in Algorithm 1. For the subset used in prompt selection during the $\underline{\mathbb{UT}}$ step, we consider training tasks in $\mathcal{D}^{(r)}_{\text{cont}}$ where the agent succeeds without compression but fails with compression. For the $\underline{\mathbb{CO}}$ step, we use training tasks in $\mathcal{D}^{(r)}_{\text{succ}}$ where the agent succeeds with compression. We perform one round consisting of a single $\underline{\mathbb{UT}}$ step and a single $\underline{\mathbb{CO}}$ step to obtain the guidelines used in our experiments, unless otherwise noted.

Baselines For FIFO, we keep last 5 interaction turns which fits in similar compression rate in average with ACON. For retrieval, we also retrieve 4 interaction turns and keep the last turn. We use OpenAI text-embedding-3-large for embedding. For LLMLingua, we set keep rate as 30% for both history and observation. For naive prompting, we use the similar prompt with Lee et al. (2025) and do some human prompt engineering to specialize each prompt to history or observation compression.

Compressor & Agent Distillation Both compressor and agent distillation use LoRA (Hu et al., 2022) with rank 16, $\alpha=32$, learning rate 10^{-4} , 3 epochs, batch size 4, and maximum sequence length 10,000. We adopt linear warmup (5% ratio), weight decay 0.01, and AdamW optimizer. No hyperparameter tuning was performed; the same setup is applied across all models and benchmarks.

https://openrouter.ai/

We sample a single generation from the teacher for fine-tuning, leaving potential improvements from hyperparameter tuning or multi-sample training for future work. We use 1 A100 80GB GPU for both training and inference. For inference of fine-tuned models, we use greedy decoding (temperature 0.0).

C ADDITIONAL RESULTS

We provide additional quantitative results to complement the main experiments in Section 4.

OfficeBench difficulty breakdowns. We further analyze OfficeBench with gpt-4.1 by difficulty level. The detailed breakdown in Table 8 shows that ACON yields the largest gains on the most challenging tasks in Level 3.

Experiments with gpt-4.1-mini. Results for the smaller variant gpt-4.1-mini (OpenAI, 2025b) across three benchmarks are reported in AppWorld (Table 9), OfficeBench (Table 10), and 8-objective QA (Table 11). The trends of ACON are consistent with those for gpt-4.1 in Section 4. In particular, Table 9 shows that history compression improves the performance of gpt-4.1-mini compared to the baseline, complementing the findings in Section 4.4 that ACON enhances the effectiveness of smaller LM agents. These results highlight the robustness of our method under resource-constrained settings.

Experiments with gpt-5-chat. We also evaluate on AppWorld using gpt-5-chat (OpenAI, 2025c), as reported in Table 12. The improvements follow the same trend as with gpt-4.1, demonstrating that ACON generalizes to the latest stronger proprietary models.

Distilled optimizer. Additional results for the distilled optimizer in AppWorld are shown in Table 13. Beyond the analysis in Section 4.3, we also include experiments where the compressor is distilled using guidelines without optimization. The results confirm that optimized guidelines consistently yield stronger performance when distilled into smaller models.

History and observation compression. In Table 14, we report ablations with gpt-4.1 using both history and observation compression. While combining the two compressions achieves larger reductions in peak token usage and dependency, it also leads to substantial performance degradation compared to applying either compression alone.

Additional guideline optimization step. We investigate whether running an extra utility maximization step (UT) after the standard sequence of utility maximization and compression maximization (CO) is beneficial. As shown in Table 14, this additional iteration results in a performance drop, indicating that a single round of optimization is sufficient for effective guideline learning.

Distilled compressor for observation. In addition to Section 4.3, we report results for observation compressor distillation in Figure 8. Similar to history compression, the performance is largely preserved after distillation, confirming that ACON enables effective transfer of optimized observation compressors to smaller models.

Case study: history compression turns failure into success. A notable case study illustrates how history compression enables a smaller agent to succeed on tasks that would otherwise fail. In the uncompressed trajectory in Example E.2, the gpt-4.1-mini agent repeatedly attempted to use the file_system APIs without managing authentication, leading to persistent 401 Unauthorized errors. After compressing the history as in Example E.3, however, the compressed history retained only the essential reasoning steps: the need for both username and password, the importance of passing the returned access_token into subsequent calls, and the absence of proxy APIs in the supervisor app.

This compressed context prevented redundant exploration and guided the agent directly to the correct sequence—login with full credentials, capture the token, and provide it explicitly in show_directory and delete_file calls. As a result, the agent was able to enumerate and

remove all .pdf files in /downloads, a task it had previously failed. This example highlights how compression does not merely shorten history but clarifies critical dependencies, turning a failure trajectory into a successful one.

Table 7: Comparison between ACON and MEM1. The two methods operate under different assumptions regarding model accessibility, training requirements, and architectural coupling.

Dimension	ACON (ours)	MEM1 (Zhou et al., 2025)
Difficusion	ACON (ours)	VIEWI (Zhou et al., 2023)
Is model training not required?	✓ no model training or weight updates required	X requires RL training on the agent model
Can the method work without access to model weights?	✓ works with open-source and proprietary API models	✗ requires full model access and gradients
Can the agent and compressor be different models?	✓ supports decoupled design with different model sizes	✗ reasoning and compression are integrated into a single model
Is it possible to use a large agent with a small compressor?	✓ supports combinations like gpt-4.1 agent + Qwen3-14B compressor	✗ same model must serve as both agent and compressor
Does optimization avoid GPU-based RL cost?	✓ under \$2 for guideline optimization, no GPU needed	★ RL policy training requires multiple trajectories and GPU computation

Comparison with MEM1. MEM1 (Zhou et al., 2025) proposes a learnable context compression policy trained jointly with the agent through reinforcement learning. This design couples reasoning and compression within a single trainable model and requires full access to model weights and gradient updates. In contrast, our method can perform optimization entirely at the prompt-level without any weight updates, enabling the agent and compressor to be different models.

This decoupling allows combinations that are not possible in MEM1. For example, one can use a large proprietary model such as gpt-4.1 as the agent while employing a smaller open-source model such as Qwen3-14B as the compressor after distillation, a configuration that MEM1 cannot support due to its unified training requirement. This flexibility makes ACON applicable to both open-source and proprietary API-based models, including settings where model internals are inaccessible. A detailed comparison is summarized in Table 7.

D QUALITATIVE EXAMPLES

We complement the quantitative results with qualitative illustrations.

Compression guidelines. We present examples of compression guidelines before and after optimization in AppWorld. The history compression guideline before optimization is shown in Prompt E.5, the optimized version (UT) in Prompt E.6, and the optimized version (UTCO) in Prompt E.7. Similarly, observation compression guideline examples are provided in Prompt E.8 and Prompt E.9, and the optimized version (UTCO) in Prompt E.10. These comparisons demonstrate that optimization yields more targeted guidelines for compressors.

Compressed histories. Compression Example E.1 illustrates history segments before and after guideline optimization in AppWorld with gpt-4.1. The optimized guideline retains a more detailed record of task progress, including variable states and guardrails for the environment. After the compression maximization step (CO), the histories become shorter while still preserving the essential information required for future decision-making. This qualitative evidence demonstrates how our framework improves both the efficiency and effectiveness of context compression, complementing the guideline optimization procedure described in Section 3.3.

We also present Compression Example E.2 for 8-objective QA and Compression Example E.3 for OfficeBench, which confirm that the effects of guideline optimization are consistent across benchmarks.

Compressed observations. Compression Example E.4 shows observations before and after guide-line optimization in AppWorld. We illustrate the case of printing available APIs for the Spotify app, which produces a lengthy observation. The optimized guideline yields a more structured and faithful representation: whereas naive prompting loses the JSON format and omits the crucial "play_music" API, the optimized version preserves both structure and key functionality necessary to complete the task

E LLM USAGE

We used large language models (LLMs) solely as a writing assistant, for improving grammar and clarity of the paper. No part of the research ideation, experimental design, or analysis relied on LLMs

```
1243
1244
1245
               Algorithm 1 Alternating Guideline Optimization (\overline{UT} \leftrightarrow \overline{CO})
1246
               Input: Training set indices \mathcal{I}; fixed agent \mathcal{M}(\cdot; \theta, \mathcal{P}_{agent}); compressor f(\cdot; \phi, \mathcal{P}); initial guideline
1247
                      \mathcal{P}^{(0)}; tradeoff \lambda \geq 0; rounds R; candidates K
1248
               Output: Optimized guideline \mathcal{P}^*
1249
                      Notation. For each i \in \mathcal{I} and guideline \mathcal{P}:
1250
                           baseline (no compression): context sequence H_i with success r_i^{\text{base}} \in \{0, 1\}
1251
                           compressed: H'_i(\mathcal{P}) with success r_i(\mathcal{P}) \in \{0,1\} and cost C(H'_i(\mathcal{P})) = \sum_t C(h'_{i,t-1},o'_{i,t})
1252
                      // 0) Collect baseline contexts (no compression)
1253
                1: for all i \in \mathcal{I} do
1254
                            Run \mathcal{M} without compression to obtain \mathbf{H}_i and r_i^{\text{base}}
1255
                3: end for
                4: \mathcal{I}^+ \leftarrow \{ i \in \mathcal{I} \mid r_i^{\mathsf{base}} = 1 \}
                                                                                                                               1256
                5: for r = 0 to 2R - 2 step 2 do
1257
                      // Stage A: \overline{UT} (reward-first update using H vs H')
                            for all i \in \mathcal{I} do
                6:
                                   Run \mathcal{M} with compression f(\cdot; \phi, \mathcal{P}^{(r)}) to obtain H'_i(\mathcal{P}^{(r)}), r_i(\mathcal{P}^{(r)}), C(H'_i(\mathcal{P}^{(r)}))
                7:
                8:
1261
                            \mathcal{D}_{\mathsf{cont}}^{(r)} \leftarrow \{\, (\boldsymbol{H}_i, \boldsymbol{H}_i'(\mathcal{P}^{(r)})) \mid i \in \mathcal{I}^+, \; r_i(\mathcal{P}^{(r)}) = 0 \,\}
                9:
1262
                            for all (\boldsymbol{H},\boldsymbol{H}')\in\mathcal{D}_{\mathsf{cont}}^{(r)} do
                                                                                                     \triangleright contrastive feedback: what did H' miss vs H?
               10:
1263
                                   Feedback \leftarrow LLM(FeedbackInstr, H, H')
               11:
1264
                                   Append to multiset \mathcal{F}_{util}
               12:
1265
                            end for \{\mathcal{P}_k^{(r+1)}\}_{k=1}^K \leftarrow \text{LLM}(\text{UpdateInstr}, \mathcal{P}^{(r)}, \parallel_{f \in \mathcal{F}_{\text{util}}} f)
               13:
1266
               14:
                                                                                                                                                          // ||: concatenation
1267
                            Select by reward: evaluate on a held-out subset of \mathcal{I}^+ and pick
               15:
1268
                                                            k_{\mathsf{util}}^* \leftarrow \arg\max_{k} \; \mathsf{SuccessRate} \big( \{ r_i(\mathcal{P}_k^{(r+1)}) \}_{i \in \mathcal{I}^+} \big)
1269
1270
                            \mathcal{P}_{\text{util}}^{(r+1)} \leftarrow \mathcal{P}_{k_{\text{util}}^*}^{(r+1)}
1271
               16:
1272
                      // Stage B: \overline{CO} (cost-minimizing refinement using only H')
1273
                            \begin{array}{l} \textbf{for all } i \in \mathcal{I} \textbf{do} \\ \textbf{Using } \mathcal{P}_{\textbf{util}}^{(r+1)}, \textbf{obtain } \boldsymbol{H}_i', r_i, C(\boldsymbol{H}_i') \end{array}
               17:
1274
               18:
                            end for
1275
               19:
               20:
                            \mathcal{D}_{\mathsf{succ}}^{(r)} \leftarrow \{ \, \boldsymbol{H}_i' \mid r_i = 1 \, \}
1276
                            for all oldsymbol{H}' \in \mathcal{D}_{	ext{succ}}^{(r)} do
1277
               21:
                                                                                                                                \triangleright find redundant spans within H'
1278
               22:
                                   CompFeedback \leftarrow LLM(CompressInstr, H')
               23:
                                   Append to multiset \mathcal{F}_{comp}
1279
               24:
                            end for
1280
                            \{\widetilde{\mathcal{P}}_k^{(r+2)}\}_{k=1}^K \leftarrow \text{LLM}(\text{UpdateInstr\_Compress}, \mathcal{P}_{\text{util}}^{(r+1)}, \parallel_{f \in \mathcal{F}_{\text{comp}}} f)
Select by reward–cost: evaluate on a held-out split of \mathcal{I} and pick
               25:
1281
1282
1283
                          k_{\mathsf{comp}}^* \leftarrow \arg\max_{k} \left( \operatorname{SuccessRate}(\{r_i(\tilde{\mathcal{P}}_k^{(r+2)})\}) - \lambda \cdot \operatorname{NormCost}(\{C(\boldsymbol{H}_i'(\tilde{\mathcal{P}}_k^{(r+2)}))\}) \right)
1284
1285
                            \mathcal{P}^{(r+2)} \leftarrow \tilde{\mathcal{P}}_{k_{\mathsf{comp}}^*}^{(r+2)}
               27:
1286
               28:
                            if early-stop criterion satisfied then
                                                                                                         ▷ e.g., success/cost convergence or budget met
1287
               29:
                                   break
1288
                            end if
               30:
1289
               31: end for
1290
               32: \mathcal{P}^* \leftarrow \mathcal{P}^{(r+2)}
1291
               33:
               34: return \mathcal{P}^*
1293
```

Table 8: Detailed results on **OfficeBench** benchmark. We report accuracy (%), and efficiency metrics: average steps, peak input tokens (10^3) , and dependency (10^6) for Average and each difficulty level. Best values are in bold. Rows in blue background indicate the results from **ours**.

		Averag	ge (All)		Leve	l 1 (1-app	o, 42)	Leve	l 2 (2-app	, 22)	Leve	d 3 (3-app	, 31)
Method	Acc. ↑	Steps ↓	Peak ↓	Dep. ↓	Acc. ↑	Peak ↓	Dep. ↓	Acc. ↑	Peak ↓	Dep. ↓	Acc. ↑	Peak ↓	Dep.↓
Agent: gpt-4.1/Compressor: gpt-4.1													
No Compression	76.84	11.52	7.27	4.43	92.86	6.23	4.05	77.27	6.14	1.81	54.84	8.37	6.08
History Compres	sion												
FIFO	67.37	12.26	4.02	2.64	83.33	4.19	0.72	63.64	3.51	1.01	48.39	4.23	4.39
Retrieval	65.26	16.20	4.33	2.06	85.71	4.35	0.84	63.64	3.52	1.37	38.71	4.78	2.99
LLMLingua	70.53	10.89	4.65	1.85	83.33	4.17	0.67	68.18	4.61	1.18	54.84	4.88	2.74
Prompting	71.58	10.13	4.40	1.10	85.71	4.18	0.81	77.27	4.53	1.08	48.39	4.42	1.23
ACON UT	74.74	13.13	4.93	3.85	85.71	4.71	6.89	72.73	4.64	1.44	61.29	5.19	3.89
ACON UTCO	72.63	11.54	4.54	1.91	88.10	3.92	0.76	72.73	4.72	1.16	51.61	4.71	2.84
Observation Con	pression												
LLMLingua	71.58	11.89	7.38	6.14	80.95	7.35	12.40	72.73	6.31	2.11	58.06	7.99	5.70
Prompting	55.79	12.24	6.44	2.68	78.57	4.51	0.98	50.00	6.98	2.61	29.03	6.98	3.46
ACON UT	73.68	10.83	6.55	3.85	90.48	6.57	8.02	77.27	6.11	1.97	48.39	6.80	3.10
ACON UTCO	72.63	10.28	6.17	2.88	88.10	4.75	0.82	72.73	6.41	2.09	51.61	6.65	4.22

Table 9: Results across different difficulty levels on **Appworld** benchmark (test-normal) with gpt-4.1-mini. We adopt the same compression guidelines as those used in the gpt-4.1 experiments. Each block reports accuracy (task goal completion score), average steps, average peak input tokens (10^3), and average dependency (10^6) for agents. Best results in each column are highlighted in bold. Rows in blue background indicate the results from **ours** .

Method		Average	e (168)			Easy (57)		M	edium (4	3)		Hard (63)	
	Acc. ↑	Steps ↓	Peak ↓	Dep.↓	Acc. ↑	Peak ↓	Dep.↓	Acc. ↑	Peak ↓	Dep.↓	Acc. ↑	Peak ↓	Dep.↓
			Agent:	gpt-4.	1-mini	/ Compre	essor: gp	ot-4.1-	mini				
No compression	35.7	18.14	8.55	5.07	56.1	6.45	3.72	31.2	8.31	4.79	20.6	10.64	9.18
History Compres	ssion												
FIFO	39.3	30.39	6.18	5.24	75.4	4.76	2.66	35.4	5.33	4.81	9.5	8.10	7.91
Retrieval	14.9	40.18	7.49	5.95	36.8	7.10	4.29	8.3	7.44	6.80	0.0	7.89	6.81
LLMLingua	36.3	28.41	7.24	6.65	66.7	6.96	3.84	33.3	7.05	7.60	11.1	7.62	8.47
Prompting	35.7	24.98	6.56	4.95	64.9	5.96	2.90	27.1	6.65	5.35	15.9	6.84	6.49
ACON <u>UT</u>	42.3	22.46	6.51	5.48	64.9	5.87	2.62	37.5	7.18	5.22	25.4	7.18	8.25
ACON UTCO	32.7	24.27	6.99	4.97	57.9	7.50	2.77	33.3	8.45	4.99	9.5	6.95	6.97
Observation Con	npression	1											
LLMLingua	25.6	20.75	8.04	8.21	38.6	6.13	3.03	27.1	8.74	13.78	12.7	9.24	8.65
Prompting	33.9	16.71	6.04	3.87	59.7	5.21	3.41	33.3	5.99	3.27	11.1	6.83	4.74
ACON UT	33.9	16.78	6.86	4.58	59.7	5.44	2.93	33.3	7.13	4.26	11.1	7.97	6.38
ACON UTCO	27.4	17.89	6.37	4.44	40.4	5.18	2.40	35.4	6.84	5.03	9.5	7.09	5.82

Table 10: Detailed results on **OfficeBench** benchmark with gpt-4.1-mini. We adopt the same compression guidelines as those used in the gpt-4.1 experiments. We report accuracy (%), and efficiency metrics: average steps, peak input tokens (10^3) , and dependency (10^6) for Average and each difficulty level. Rows in blue background indicate the results from **ours** .

		Averag	ge (All)		Leve	el 1 (1-app	o, 42)	Leve	l 2 (2-app	o, 22)	Leve	l 3 (3-app	, 31)
Method	Acc. ↑	Steps ↓	Peak ↓	Dep. ↓	Acc. ↑	Peak ↓	Dep. ↓	Acc. ↑	Peak ↓	Dep. ↓	Acc. ↑	Peak ↓	Dep. ↓
			Agent	:gpt-4	.1-mini	/ Compr	essor: gp	t-4.1-r	mini				
No Compression	72.63	11.96	7.36	3.92	88.10	6.66	4.29	68.18	4.97	1.01	54.84	9.02	5.40
History Compres	sion												
FIFO	65.26	10.91	4.03	1.46	83.33	4.10	0.78	59.09	3.69	0.96	45.16	4.19	2.03
Retrieval	67.37	14.46	4.55	2.74	85.71	5.85	5.86	59.09	3.47	0.87	48.39	4.59	2.45
LLMLingua	67.39	11.59	4.90	2.18	87.18	4.31	3.87	59.09	4.58	0.92	48.39	5.34	2.17
Prompting	71.58	11.78	4.93	3.10	85.71	4.73	4.75	72.73	4.40	0.86	51.61	5.32	3.06
ACON	73.68	12.41	4.82	1.96	88.10	4.12	0.83	68.18	4.39	0.86	58.06	5.37	3.07
Observation Com	pression												
LLMLingua	66.32	11.02	6.34	2.40	78.57	6.09	2.12	63.64	4.82	0.97	51.61	7.30	3.34
Prompting	73.68	11.43	6.45	2.62	88.10	4.82	1.44	72.73	4.95	1.06	54.84	8.01	4.01
ACON	71.58	10.96	6.00	2.19	88.10	4.45	1.06	63.64	4.89	1.00	54.84	7.30	3.36

Table 11: Results on **8-objective QA** benchmark with gpt-4.1-mini. We adopt the same compression guidelines as those used in the gpt-4.1 experiments. We report EM/F1 and efficiency metrics (Steps, Peak input tokens (10^3) , and Dependency (10^6)).

Method	EM ↑	F1 ↑	Steps ↓	Peak ↓	Dep. ↓
Agent: gpt-4	l.1-min	i / Com	pressor: g	pt-4.1-	mini
No compression	0.330	0.436	19.80	12.93	5.63
History Compre	ssion				
FIFO	0.024	0.031	28.45	5.33	3.89
Retrieval	0.143	0.190	26.90	5.34	3.55
LLMLingua	0.140	0.194	25.24	6.69	3.92
Prompting	0.149	0.207	25.27	4.85	2.44
ACON	0.238	0.325	21.05	4.78	2.03
ACON (iter2)	0.248	0.353	19.18	4.79	1.79
Observation Cor	npression	1			
LLMLingua	0.316	0.430	15.96	5.54	1.60
Prompting	0.282	0.402	11.71	3.91	0.65
ACON	0.323	0.434	14.42	4.71	1.10
ACON (iter2)	0.316	0.443	11.69	3.97	0.63

Table 12: Results across different difficulty levels on **AppWorld** benchmark (test-normal) with gpt-5-chat. We adopt the same compression guidelines as those used in the gpt-4.1 experiments. Each block reports accuracy (task goal completion score), steps, peak input tokens (10^3) , and dependency (10^6) for agents. Best results in each column are highlighted in bold. Rows in blue background indicate the results from **ours** .

Method		Average	(168)			Easy (57)		M	edium (48	3)		Hard (63)	
	Acc. ↑	Steps ↓	Peak ↓	Dep.↓	Acc. ↑	Peak ↓	Dep.↓	Acc. ↑	Peak ↓	Dep.↓	Acc. ↑	Peak ↓	Dep.↓
			Age	nt: gpt-	-5-chat	/ Compr	essor: gp	t-5-ch	at				
No compression	66.7	16.45	9.67	4.78	89.5	7.55	2.31	64.6	9.58	4.13	47.6	11.67	7.51
History Compres	ssion												
FIFO (last-5)	46.4	30.61	6.81	4.85	79.0	5.21	2.10	43.8	6.82	5.50	19.1	8.24	6.84
Prompting	58.9	22.24	7.46	4.02	82.5	7.15	2.13	66.7	7.19	3.69	31.8	7.93	5.97
ACON UT	58.3	20.15	6.97	3.74	80.7	6.66	2.04	66.7	7.08	3.40	31.8	7.16	5.54
ACON <u>UTCO</u>	62.5	22.29	7.26	3.85	86.0	6.44	2.04	72.9	6.98	3.93	33.3	8.20	5.42
Observation Con	npression	1											
Prompting	60.1	17.39	6.50	3.72	80.7	4.98	1.72	68.8	6.40	3.48	34.9	7.96	5.70
ACON UT	65.5	17.16	7.58	3.96	84.2	5.62	1.94	68.8	7.49	3.46	46.0	9.41	6.16
ACON <u>UTCO</u>	62.5	18.21	7.21	4.24	80.7	5.52	2.02	70.8	7.18	3.69	39.7	8.76	6.67
History + Observ	ation Co	mpression	1										
ACON UT	63.1	20.02	5.89	3.63	77.2	5.27	1.92	77.1	6.03	3.52	39.7	6.35	5.28
ACON UTCO	58.9	22.90	5.83	4.07	80.7	5.35	1.94	77.1	5.94	3.56	25.4	6.17	6.39

Table 13: Results across different difficulty levels on **AppWorld** with **distilled compressors**. We report accuracy (task goal completion score), average steps, peak input tokens (10^3), and dependency (10^6). For all compressors, we use the optimized compression guideline after the utilization maximization $\overline{\text{UT}}$ step. 'Fine-tune' means that we fine-tune small models with outputs from naive prompt before compression guideline optimization.

Method	Average				Easy			Medium			Hard		
Memod	Acc. ↑	Steps ↓	Peak ↓	Dep.↓	Acc. ↑	Peak ↓	Dep.↓	Acc. ↑	Peak ↓	Dep.↓	Acc. ↑	Peak ↓	Dep.↓
Agent: gpt-4.1/Compressor: gpt-4.1-mini or Distilled models (Qwen3, Phi-4)													
History Compression													
Prompting (gpt-4.1-mini)	39.3	23.61	7.03	5.19	64.9	6.64	3.17	35.4	7.63	5.42	19.1	6.93	6.84
ACON (gpt-4.1-mini)	47.6	21.46	7.25	5.24	75.4	6.75	2.84	35.4	7.25	5.36	31.8	7.70	7.32
Fine-tune (Qwen3-14B)	44.6	24.16	7.16	4.95	71.9	6.79	2.88	43.8	7.39	4.88	20.6	7.33	6.88
ACON (Qwen3-14B)	50.0	21.72	6.83	4.80	79.0	6.42	2.54	50.0	6.87	4.89	23.8	7.17	6.79
ACON (Qwen3-8B)	47.0	21.58	6.98	4.76	71.9	6.64	2.93	37.5	7.24	4.67	31.8	7.09	6.48
ACON (Phi-4)	44.6	21.19	7.24	4.76	68.4	7.33	2.75	39.6	7.12	4.16	27.0	7.26	7.04
Observation Compression													
Prompting (gpt-4.1-mini)	44.0	16.67	6.84	4.30	71.9	5.08	2.19	35.4	6.72	3.77	25.4	8.53	6.61
ACON (gpt-4.1-mini)	48.2	18.00	8.66	6.62	71.9	6.05	2.60	37.5	9.23	7.41	34.9	10.60	9.65
Fine-tune (Qwen3-14B)	40.5	17.71	6.64	4.38	64.9	4.91	1.97	31.2	6.72	4.05	25.4	8.16	6.81
ACON (Qwen3-14B)	56.5	16.78	7.57	5.06	82.5	5.69	2.20	54.2	7.39	4.46	34.9	9.40	8.10
ACON (Qwen3-8B)	48.2	16.10	7.33	4.82	71.9	5.49	2.03	50.0	7.20	4.20	25.4	9.10	7.82
ACON (Phi-4)	50.6	16.88	7.88	5.41	77.2	5.85	2.88	52.1	7.75	4.77	25.4	9.83	8.18

Table 14: Additional results for additional guideline optimization step and unified compression on **Appworld** benchmark (test-normal). Each block reports accuracy (task goal completion score), steps, peak input tokens (10^3), and dependency (10^6) for agents. Best results in each column are highlighted in bold. Rows in blue background indicate the results from **ours** .

Method	Average (168)				Easy (57)			Medium (48)			Hard (63)		
	Acc. ↑	Steps ↓	Peak ↓	Dep.↓	Acc. ↑	Peak ↓	Dep.↓	Acc. ↑	Peak ↓	Dep.↓	Acc. ↑	Peak ↓	Dep.↓
Agent: gpt-4.1/Compressor: gpt-4.1													
History Compression													
ACON UTCOUT	47.0	22.28	7.22	4.66	68.4	7.01	2.69	58.3	7.16	4.39	19.1	7.45	6.65
History + Observation Compression													
Prompting	36.3	19.33	5.38	3.44	71.9	4.87	1.80	21.6	5.63	3.60	14.3	5.64	4.79
ACON	45.8	20.32	5.85	4.26	75.4	5.29	2.07	39.6	6.15	4.29	23.8	6.12	6.21
ACON UTCO	44.6	21.75	5.90	4.98	77.2	5.50	2.33	39.6	6.18	3.80	19.1	6.18	8.28

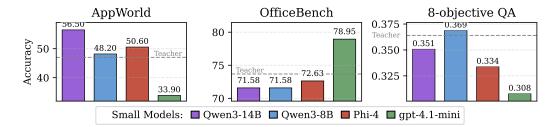


Figure 8: **Results of distilled compressors on observation compression** with gpt-4.1 as the agent. Student models (Qwen3-14B, Qwen3-8B, Phi-4) are distilled from gpt-4.1 compressor using the optimized compression guideline after <u>UT</u> step, and evaluated across all benchmarks. We also include result with gpt-4.1-mini without distillation for comparison.

```
1458
         Prompt E.1: Prompt for analysis before prompt optimization (utility step)
1459
1460
         You are an expert agent trajectory auditor.
1461
1462
         Analyze why the HISTORY-OPTIMIZED agent failed OR became
1463
             significantly less efficient while the BASELINE succeeded.
1464
         You are given:
         - task_name: {{ task_name }}
1465
         - Baseline full history (single continuous session)
1466
         - Optimized history split into multiple sessions where each new
1467
             session starts with a fresh system + user prompt and an
1468
             injected <HISTORY_SUMMARY> summarizing earlier interactions.
1469
         - baseline_success={{ baseline_success }} optimized_success={{
1470
             optimized_success }}
         - baseline_env_steps={{ baseline_env_steps | default('null') }}
1471
             optimized_env_steps={{ optimized_env_steps | default('null')}
1472
              }} step_ratio={{ step_ratio | default('null') }}
1473
             performance_regression={{ performance_regression | default('
1474
             false') }}
1475
         Goals:
1476
         1. Determine whether summarization / session resetting removed,
1477
             distorted, delayed, or bloated reasoning causing failure OR
1478
             inflated step count (> threshold factor of baseline).
1479
         2. Identify the FIRST divergence point where the optimized
             trajectory meaningfully deviates from the successful &
1480
             efficient baseline path.
1481
         3. Categorize root causes (e.g., Missing Critical Fact,
1482
             Incorrect Summary, Lost Variable/State, Unnecessary Re-
1483
             discovery, Instruction Drift, API Misuse, Premature
1484
             Completion, Token Truncation, Inefficient Looping, Redundant
              API Calls, Over-Exploration, Other).
1485
         4. Extract concrete evidence snippets (quote exact lines) from
1486
             baseline vs optimized showing:
1487
            - Critical facts present in baseline but absent/altered in
1488
             optimized (esp. after a session boundary)
1489
            - Summary inaccuracies (baseline ground truth vs summary text
             )
1490
            - Redundant or looping action patterns causing step inflation
1491
1492
         5. Suggest precise remediation strategies: summary style changes
1493
             , retain variable/value tables, move session boundaries,
1494
             guardrail prompts, caching, early-exit heuristics, loop
             detection, etc.
1495
         6. Provide a reliability_score (0.0-1.0) reflecting confidence
1496
             in your causal attribution.
1497
         7. If performance_regression == true, analyze efficiency
1498
             degradation even if optimized_success==true.
1499
         Output STRICTLY valid JSON object with keys:
1500
1501
            "task_name": str,
1502
           "divergence_step_description": str,
1503
            "root_cause_categories": [str, ...],
1504
            "missing_or_distorted_facts": [ {"baseline": str, "
             optimized_context_absent_or_changed": str, "impact": str} ],
1505
            "summary_inaccuracies": [ {"summary_excerpt": str, "issue_type
1506
             ": str, "correct_baseline_reference": str, "impact": str} ],
1507
            "lost_state_variables": [ {"name_or_pattern": str,
             baseline_evidence": str, "optimized_issue": str} ],
            "api_or_action_errors": [ {"optimized_step_excerpt": str, "
1509
             error_type": str, "improvement": str} ],
1510
```

```
1512
            "inefficiency_patterns": [ {"pattern": str, "evidence_excerpt
1513
             ": str, "excess_steps": int, "cause": str, "remediation":
1514
             str}],
1515
            "timeline_of_divergence": [ {"phase": str, "optimized_excerpt
1516
             ": str, "baseline_contrast": str, "effect": str}],
1517
            "performance regression": bool,
            "baseline_env_steps": int | null,
1518
            "optimized_env_steps": int | null,
1519
            "step_ratio": float | null,
1520
            "remediation_recommendations": [ str, ... ],
1521
            "recovery_opportunities_missed": [ { "optimized_excerpt": str,
1522
             "missed_fix_action": str} ],
            "reliability_score": float,
1523
            "concise_failure_mechanism_summary": str
1524
1525
1526
         If some sections have no data, use an empty list. For non-
             applicable numeric fields use null.
         Do NOT include any extra commentary outside JSON.
1528
1529
1530
         BASELINE_HISTORY_START
1531
         {{ baseline_history }}
1532
         BASELINE_HISTORY_END
1533
         OPTIMIZED_MULTI_SESSION_HISTORY_START
1534
         {{ optimized_history }}
1535
         OPTIMIZED_MULTI_SESSION_HISTORY_END
1536
1537
         Failure or performance report / metadata (may be null):
1538
         {{ failure_report }}
1539
         Proceed with rigorous comparison.
1540
1541
```

Prompt E.2: Prompt for prompt optimization after analysis (utility step)

1542

```
1544
         You are an expert prompt engineer tasked with refining a HISTORY
1545
              SUMMARIZATION prompt.
1546
         Rewrite the ORIGINAL PROMPT to reduce length of the HISTORY
1547
             SUMMARY while preserving factual continuity for the next
1548
             session.
         Ground all changes in the PER-SAMPLE REDUCTION SIGNALS below. Do
1549
              not aggregate across samples; use the patterns and rules as
1550
             -is.
1551
1552
         Constraints:
1553
         - Keep all Jinja placeholders, variable names, and structure
             intact where possible.
1554
         - Add explicit, concrete rules that prevent verbosity and retain
1555
              essential state.
1556
         - Do not include literal values from prior content; refer to
1557
             variable names only.
1558
         - Output ONLY the improved prompt template (no extra commentary)
1559
1560
         Context (samples below are the only ground truth signals to use)
1561
1562
         - Average original summary size (chars) across sampled set: {{
1563
             avg_orig_chars }}
1564
         {% for s in samples %}
1565
```

1598 1599

1601

1602

1603

1604

1605

1606

1607

1608

1609

1610

1611

1612

1613 1614 1615

1616 1617

1618 1619

```
1566
         ==== SAMPLE {{ loop.index0 }} =====
1567
         - Task/Session: {{ s.task_label }} / {{ s.session or 'unknown-
1568
             session' }}
1569
         - Analysis Overview:
1570
         {% if s.overview %}
1571
         {% for k, v in s.overview.items() %} - {{ k }}: {{ v }}
1572
         {% endfor %}
         {% else %} - (none provided)
1573
         {% endif %}
1574
1575
         - Removals (patterns -> action):
1576
         {% for r in s.removals %} - [{{ r.category | default('unknown')}
              }}] {{ r.pattern | default('') }} -> {{ r.action | default
1577
             ('drop') }}
1578
         {% endfor %}
1579
1580
         - KEEP examples (evidence-driven essentials):
1581
         % = { (k.reason | default('') ) } 
              Evidence: {{ k.evidence_spans | default([]) | join('; ') }}
1582
         {% endfor %}
1583
1584
         - Summary Rules:
1585
         {% for rule in s.rules %} - {{ rule }}
1586
         {% endfor %}
1587
         {% endfor %}
1588
1589
         Original Prompt Template (verbatim between markers):
1590
         <<<ORIGINAL_PROMPT>>>
1591
         {{ original_prompt }}
         <<<ORIGINAL_PROMPT>>>
1592
1593
         Output only the improved prompt template text, ready to be used
1594
             as a Jinja template.
1595
1596
```

Prompt E.3: Prompt for analysis before prompt optimization (compression step)

```
You are an expert prompt engineer tasked with refining a HISTORY
    SUMMARIZATION prompt.
Rewrite the ORIGINAL PROMPT to reduce length of the HISTORY
   SUMMARY while preserving factual continuity for the next
   session.
Ground all changes in the PER-SAMPLE REDUCTION SIGNALS below. Do
    not aggregate across samples; use the patterns and rules as
   -is.
Constraints:
- Keep all Jinja placeholders, variable names, and structure
   intact where possible.
- Add explicit, concrete rules that prevent verbosity and retain
    essential state.
- Do not include literal values from prior content; refer to
   variable names only.
- Output ONLY the improved prompt template (no extra commentary)
Context (samples below are the only ground truth signals to use)
- Average original summary size (chars) across sampled set: {{
   avg_orig_chars }}
```

```
1620
          {% for s in samples %}
1621
          ==== SAMPLE {{ loop.index0 }} =====
1622
          - Task/Session: {{ s.task_label }} / {{ s.session or 'unknown-
1623
              session' }}
1624
          - Analysis Overview:
1625
          {% if s.overview %}
          {% for k, v in s.overview.items() %} - {{ k }}: {{ v }}
1626
          {% endfor %}
1627
          {% else %} - (none provided)
1628
          {% endif %}
1629
1630
          - Removals (patterns -> action):
          {% for r in s.removals %} - [{{ r.category | default('unknown')}
1631
               }}] {{ r.pattern | default('') }} -> {{ r.action | default
1632
              ('drop') }}
1633
          {% endfor %}
1634
1635
          - KEEP examples (evidence-driven essentials):
          \{\% \text{ for } k \text{ in s.keeps } \%\} - Reason: \{\{ k.reason \mid default('') \}\} \mid \emptyset\}
1636
               Evidence: {{ k.evidence_spans | default([]) | join('; ') }}
1637
          {% endfor %}
1638
1639
          - Summary Rules:
1640
          {% for rule in s.rules %} - {{ rule }}
          {% endfor %}
1641
1642
          {% endfor %}
1643
1644
          Original Prompt Template (verbatim between markers):
1645
          <<<ORIGINAL_PROMPT>>>
          {{ original_prompt }}
1646
          <<<ORIGINAL_PROMPT>>>
1647
1648
          Output only the improved prompt template text, ready to be used
1649
              as a Jinja template.
1650
1651
```

Prompt E.4: Prompt for analysis before prompt optimization (compression step)

1652

1653 1654

1655

1656

1657

1658

1659

1660

1661

1662

1663

1664

1665

1666

1667

1668 1669 1670

1671

1672

```
You are an expert prompt engineer tasked with refining a HISTORY
    SUMMARIZATION prompt.
Rewrite the ORIGINAL PROMPT to reduce length of the HISTORY
   SUMMARY while preserving factual continuity for the next
   session.
Ground all changes in the PER-SAMPLE REDUCTION SIGNALS below. Do
    not aggregate across samples; use the patterns and rules as
   -is.
Constraints:
- Keep all Jinja placeholders, variable names, and structure
   intact where possible.
- Add explicit, concrete rules that prevent verbosity and retain
    essential state.
- Do not include literal values from prior content; refer to
   variable names only.
- Output ONLY the improved prompt template (no extra commentary)
Context (samples below are the only ground truth signals to use)
- Average original summary size (chars) across sampled set: {{
   avg_orig_chars }}
```

```
1674
1675
          {% for s in samples %}
1676
         ===== SAMPLE {{ loop.index0 }} =====
1677
         - Task/Session: {{ s.task_label }} / {{ s.session or 'unknown-
1678
             session' }}
1679
         - Analysis Overview:
         {% if s.overview %}
1680
          {% for k, v in s.overview.items() %} - {{ k }}: {{ v }}
1681
          {% endfor %}
1682
          {% else %} - (none provided)
1683
          {% endif %}
1684
1685
         - Removals (patterns -> action):
         {% for r in s.removals %} - [{{ r.category | default('unknown')}
1686
              }}] {{ r.pattern | default('') }} -> {{ r.action | default
1687
              ('drop') }}
1688
          {% endfor %}
1689
         - KEEP examples (evidence-driven essentials):
1690
         {% for k in s.keeps %} - Reason: {{ k.reason | default('') }} |
1691
              Evidence: {{ k.evidence_spans | default([]) | join('; ') }}
1692
          {% endfor %}
1693
1694
         - Summary Rules:
         {% for rule in s.rules %} - {{ rule }}
1695
         {% endfor %}
1696
1697
         {% endfor %}
1698
1699
         Original Prompt Template (verbatim between markers):
         <<<ORIGINAL_PROMPT>>>
1700
         {{ original_prompt }}
1701
         <<<ORIGINAL_PROMPT>>>
1702
1703
         Output only the improved prompt template text, ready to be used
1704
             as a Jinja template.
1705
```

Prompt E.5: AppWorld Prompt for history compression before optimization

```
You are maintaining a structured context-aware summary for a productivity agent. You will be given the user instruction for the agent, a list of interactions corresponding to actions taken by the agent, and the most recent previous summary if one exists. Produce the following:

### REASONING
Summarize key progress, decisions made, important observed outcomes, and rationale behind actions taken so far. Include how earlier steps influenced later ones and why certain data is retained in the summary.

### COMPLETED
List completed subtasks or successful outcomes, with brief results if applicable.

---

### [Information Source]

### USER INSTRUCTION
```

```
1728
          {{ task }}
1729
1730
          ## [PREVIOUS SUMMARY] (if any)
1731
1732
          {{ prev_summary }}
1733
          ## [HISTORY OF INTERACTIONS]
1734
1735
          {{ history }}
1736
1737
1738
         ## PRIORITIZE
1739
1740
         1. Keep all sections relevant and concise.
1741
         2. Use reusable structured formats when summarizing artifacts.
1742
         3. Ensure agent can resume task with no loss of information.
1743
         4. Include key info from errors or failed attempts to prevent
             repeated mistakes.
1744
          5. Preserve all essential artifacts and data needed to complete
1745
             the task.
1746
1747
1748
          ### [Output Format]
1749
1750
         Do **not** include the input or any additional explanation. Only
1751
              return the formatted summary.
1752
1753
1754
         Prompt E.6: AppWorld Prompt for history compression after optimization (UT)
1755
1756
         You maintain a compact, state-preserving HISTORY_SUMMARY for a
1757
             multi-session agent.
1758
1759
         Input:
          [USER INSTRUCTION] {{ task }}
1760
          [PREVIOUS SUMMARY] {{ prev_summary }}
1761
          [HISTORY OF INTERACTIONS] {{ history }}
1762
1763
         Create the following sections-use the exact headings and order:
1764
         <history_summary>
1765
1766
         1. REASONING
1767
             - Key progress, decisions, outcomes, and their rationale.
1768
             - Note how earlier steps influence later ones.
1769
         2. VARS
1770
             | name | value | purpose |
1771
             |-----|
1772
             Record every runtime value the next session must re-declare (
1773
             tokens, ids, lists, last page_index/page_limit, etc.).
1774
         3. TODO
1775
            List pending actions with enough detail to execute directly.
1776
1777
          4. COMPLETED
1778
             Bullet list of finished subtasks with brief results.
1779
          5. GUARDRAILS
1780
```

Short reminders that prevent repeat errors, e.g.

```
1782
             - Memory resets; re-create VARS before use.
1783
             - Paginate until empty page.
1784
             - Validate API parameters against spec.
1785
             - Avoid redundant logins or doc look-ups.
1786
1787
         Requirements:
          - Be concise-bullets and tables preferred; no extraneous prose.
1788
          - Preserve all essential facts, parameters, and artifacts; omit
1789
             nothing critical.
1790
           Include errors only if they inform future avoidance.
1791
          - Do not output the input or any commentary-return only <
1792
             HISTORY_SUMMARY>.
1793
1794
          Prompt E.7: AppWorld Prompt for history compression after optimization
1795
1796
         (UTCO)
1797
1798
         You maintain a compact, state-preserving HISTORY_SUMMARY for a
1799
             multi-session agent.
1800
          Input:
1801
          [USER INSTRUCTION] {{ task }}
1802
          [PREVIOUS SUMMARY] {{ prev_summary }}
1803
          [HISTORY OF INTERACTIONS] {{ history }}
1804
         Summary Compression Rules:
1805
          - Collapse multi-bullet narratives into <=2 concise sentences.
1806
          - Replace repetitive step logs with one summarizing phrase.
1807
          - Truncate long token/credential strings to "<token>" unless
1808
             verbatim reuse is required.
          - Remove unused/expired credentials, page_index/page_limit,
1809
             verbose API dumps, and table borders.
1810
          - Shrink GUARDRAILS to one bullet unless multiple items are
1811
             still critical.
1812
          - Delete tool/API log output, greetings, meta prose, and section
1813
              headers that no longer contain content.
          - Keep only variables actively referenced in upcoming steps;
1814
             list each once in VARS.
1815
          - Reference removal categories [repetition], [tool-logs], [meta
1816
             ], [formatting] to prune similar lines.
1817
          - Preserve factual continuity; never invent or alter state
1818
             variables.
          - Target summaries well under {{ max_chars | default(1500) }}
1819
             characters.
1820
1821
         Critical Essentials:
1822
         Always keep evidence-driven items required next session (e.g.,
1823
             tokens, ids, emails, amounts, lists, paths, description
             strings, brief task status).
1824
1825
         Output EXACTLY the following structure---nothing more:
1826
1827
          <history_summary>
1828
         1. REASONING
1829
            One brief paragraph on key progress and rationale.
1830
1831
1832
            key=value pairs, comma-separated; only still-needed runtime
1833
             values.
1834
          3. TODO
```

```
1836
             Bulleted next actions (<=5).
1837
1838
          4. COMPLETED
1839
             Bulleted finished subtasks (<=5).
1840
1841
          5. GUARDRAILS
             Single concise bullet, or omit if none.
1842
1843
          Return only the <HISTORY_SUMMARY> block---no additional
1844
              commentary or input echoes.
1845
1846
1847
1848
          Prompt E.8: AppWorld Prompt for observation compression before optimization
1849
1850
          Your task is to generate a "Reasoning" and a "Refined
1851
```

Observation" based on the inputs below.

In the "Reasoning", analyze the user instruction and history to identify what information from the current observation is necessary to complete the remaining steps.

Think about what parts can be summarized or transformed to reduce length, while ensuring that future actions can still be executed based on the refined observation alone.

In the "Refined Observation", include only the information that is minimal but sufficient for the next steps.

```
[Information source]
# User Instruction
{{ task }}
# History of interactions
{{ history }}
# Observation at the current time step
{{ observation }}
[Output format]
# Reasoning
... your reasoning for what matters and how to optimize it ...
# Refined Observation
... reduced and actionable observation ...
```

Prompt E.9: AppWorld Prompt for observation compression after optimization

Your task: write two sections --- "Reasoning" and "Refined Observation".

1. Reasoning

1852

1853

1854

1855

1856

1857 1858

1859

1860 1861

1862

1863 1864

1865

1866

1867

1868 1869

1870

1871

1872

1877 1878 1879

1880

1881 1882

1883

1884

1885

1886

1887

1888

- Examine task, history, and observation.
- Decide exactly which parts of the observation must be kept so the next agent step can succeed.
- Note any need to paginate (page_limit default = 5, page_index).
- Justify any data you drop.
- 2. Refined Observation

```
1890
             - Contain only the minimal yet sufficient info for the next
1891
             step.
1892
             - Always preserve:
1893
               - Every endpoint that may be called, plus its full
1894
             parameter list and defaults (especially page_limit/
1895
             page index, auth tokens).
               - Response-schema fields referenced or likely needed later
1896
              (e.g., play_count, release_date, like_count, position, ids).
               - Raw data rows required for future comparisons or loops;
1898
             if summarising, keep at least all positive-match examples.
1899
             - Never:
1900
               - Omit defaults that affect behaviour.
               - Declare parameters "not critical" without proof.
1901
               - Hallucinate endpoints or fields.
1902
               - Replace machine-readable data with vague prose.
1903
1904
          [Information source]
1905
          # User Instruction
          {{ task }}
1906
1907
          # History of interactions
1908
          {{ history }}
1909
1910
          # Observation at the current time step
          {{ observation }}
1911
1912
          [Output format]
1913
          # Reasoning
1914
          ...concise analysis explaining what is kept/removed...
1915
          # Refined Observation
1916
          ...trimmed yet complete observation ensuring future steps remain
1917
              possible...
1918
1919
1920
          Prompt E.10: AppWorld Prompt for observation compression after optimization
1921
         (UTCO)
1922
1923
          Your task: create two sections --- "Reasoning" and "Refined
1924
             Observation".
1925
1926
         1. Reasoning (<=40 words)</pre>

    Briefly state what was kept and why; note dropped

1927
             categories and any pagination needs.
1928
1929
          2. Refined Observation (use ONLY the current observation)
1930
             - Keep strictly necessary data for the next step.
1931
             - Always preserve, when present:
               - access_token or other auth values
1932
               - page_limit, page_index and other defaults that affect
1933
             calls
```

identifiers/fields needed for comparisons, loops, or API calls (ids, titles, counts, paths, etc.)
Minimise length:
Delete unused fields and narrative text.

- every endpoint name; include parameters only if required;

- Compress JSON/arrays (one object per line, no extra spaces).
- Summarise long uniform lists with a range/pattern when individual rows are not needed.
- Never invent, alter, or omit a required literal.

description <= 4 words

1934

1935

1936

1937

1938

1939

1940

1941

1942

```
1944
              - Exclude history, prior summaries, and meta comments.
1945
1946
           [Information source]
1947
           # User Instruction
1948
           {{ task }}
1949
           # History of interactions
1950
           {{ history }}
1951
1952
          \ensuremath{\text{\#}} Observation at the current time step
1953
           {{ observation }}
1954
           [Output format]
1955
          # Reasoning
1956
1957
1958
           # Refined Observation
1959
1960
1961
1962
```

```
1998
        Example E.1: AppWorld task trajectory example
1999
2000
        Problem: Name the artist least recommended to me on Spotify.
2001
2002
        # First, let's check what APIs are available under the Spotify app to
2003

ightharpoonup see if there is an API related to recommendations or artists.
        print(apis.api_docs.show_api_descriptions(app_name='spotify'))
2004
           [{"name":"login", ...}, {"name":"show_recommendations", "descriptio
2006
           → n": "Get personalized song recommendations for the user."},
2007
              . . . ]
2008
2009
        # There is an API called "show_recommendations" which provides
2010
        → personalized song recommendations. Let's check its detailed
2011
        → specification...
2012
        print(apis.api_docs.show_api_doc(app_name='spotify',
2013
        → api_name='show_recommendations'))
2014
2015
           {"path":"/recommendations", "method": "GET", "parameters": [{"name": |
2016
           → "access_token", "required":true},
           2017
2018
        # To call `show_recommendations`, I need an access token. Check the
2019
        → login spec.
2020
        print(apis.api_docs.show_api_doc(app_name='spotify',
2021
        → api_name='login'))
2022
2023
           {"path":"/auth/token", "method": "POST", "parameters": [{"name": "use |
2024

    rname"}, {"name": "password"}],
2025
           2026
2027
        # I need Timothy's Spotify account password to log in. Let's get it
2028
        \hookrightarrow from supervisor.
        spotify_passwords = apis.supervisor.show_account_passwords()
2029
        spotify_password = next((acc["password"] for acc in spotify_passwords
2030

    if acc["account_name"] == "spotify"), None)

2031
        print(spotify_password)
2032
2033
           <REDACTED>
2034
2035
        # Now log in to get the access token.
2036
        spotify_login = apis.spotify.login(username="timothy.whit@gmail.com",
2037
                                            password="<REDACTED>")
2038
        access_token = spotify_login.get("access_token")
        print(access_token)
2039
2040
           eyJhbGciOiJIUzI1NiIsInR5cCI6IkpXVCJ9...<REDACTED>...
2041
2042
2043
        # Now get the first page of personalized song recommendations.
        recommendations_page_0 = apis.spotify.show_recommendations(
2044
            access_token=access_token, page_index=0, page_limit=20
2045
2046
        print(recommendations_page_0[:4]) # show a few items only
2047
2048
```

```
2052
2053
2054
             {"song_id":10, "title": "The Curse of Loving
2055
             → You", "artists": [{"name": "Lucas Grey"}]},
             {"song_id":11,"title":"Lost in a Moment's
2056
            Grace", "artists":[{"name": "Apollo
Serenade"}, {"name": "Phoenix Rivers"}]},
2057
2058
             {"song_id":40, "title": "Burning Like a Distant

    Star", "artists":[{"name": "Aria Sterling"}]},
2060
            . . .
2061
           ]
2062
2063
        # To find the artist least recommended to Timothy, I will count how
2064
         → many times each artist appears in the recommendations. The artist
            with the fewest appearances is the least recommended.
2065
        from collections import Counter
2066
2067
        artist_names = [a["name"] for r in recommendations_page_0 for a in
2068

    r["artists"]]

2069
        artist_counts = Counter(artist_names)
        least_recommended_artist = min(artist_counts, key=artist_counts.get)
2070
        print(least_recommended_artist)
2071
2072
           Evelyn Rose
2073
2074
        # Task complete. I will submit "Evelyn Rose" as the answer.
2075
        apis.supervisor.complete_task(answer="Evelyn Rose")
2076
2077
        Example E.2: AppWorld trajectory example - gpt-4.1-mini agent fails without history compres-
2078
        sion
2079
2080
        Problem: Delete all .pdf files from ~/downloads.
2081
2082
2083
         # Understand available APIs in the file_system app.
2084
        print(apis.api_docs.show_api_descriptions(app_name='file_system'))
2085
2086
2087
           {"name": "show_account", "description": "Show your account
            → information. Unlike show_profile, this includes private
2089

    information."},
2090
           {"name": "signup", "description": "Sign up to create account."},
           { "name": "delete_account", "description": "Delete your account." },
2091
           {"name": "update_account_name", "description": "Update your first or
2092
            → last name in the account profile."},
2093
2094
           { "name": "copy_directory", "description": "Copy a directory to
2095

    another location."
},
2096
           {"name": "move_directory", "description": "Move a directory to
            → another location."},
2097
           { "name": "compress_directory", "description": "Compress a
2098

    directory."},
2099
           {"name": "decompress_file", "description": "Decompress a compressed
2100
               file."}
2101
           1
2102
2103
2104
         # List files in ~/downloads to find PDFs.
2105
```

```
2106
        downloads_path = '~/downloads'
2107
        downloads_files = apis.file_system.show_directory(path=downloads_path)
2108
        print(downloads_files)
2109
2110
           Execution failed. Traceback:
2111
2112
           Exception: Response status code is 401:
2113
           {"message": "You are either not authorized to access this
2114

→ file_system API endpoint or your access token is missing,

            → invalid or expired."}
2115
2116
2117
        # Need authentication. Check supervisor for stored credentials.
2118
2119
        print(apis.supervisor.show_account_passwords())
2120
2121
2122
           { "account_name": "amazon", "password": "Gt$!_*W"},
2123
           { "account_name": "file_system", "password": "8nNw! jZ"},
           {"account_name": "gmail", "password": "qu4Y7}s"},
2124
           {"account_name": "phone", "password": "C4n&I40"},
2125
           { "account_name": "simple_note", "password": "e+QwbmV"},
2126
           { "account_name": "splitwise", "password": "mSqG}QU" },
2127
           {"account_name": "spotify", "password": "%CCv18v"}, {"account_name": "todoist", "password": "jHZ#RPM"},
2128
           {"account_name": "venmo", "password": "Wq8!RAU"}
2129
2130
2131
2132
         # Try logging in with just the password (no username).
2133
2134
        file_system_password = [a["password"] for a in
2135
         → apis.supervisor.show_account_passwords()
        if a["account_name"] == "file_system"][0]
2136
        file_system_login_response =
2137
        → apis.file_system.login(password=file_system_password)
2138
        print(file_system_login_response)
2139
2140
           Execution failed. Traceback:
2141
2142
           Exception: Response status code is 422:
2143
           {"message":"Validation error. Reason: \nusername: Field
            → required"}
2144
2145
2146
        # Need username. Check supervisor profile for email/username.
2147
2148
        print (apis.supervisor.show_profile())
2149
2150
           {"first_name": "Susan", "last_name": "Burton", "email": "[susanmiller]
2151
            ⇔ @gmail.com] (mailto:susanmiller@gmail.com) ",
2152
           "phone_number": "3296062648", "birthday": "1994-04-30", "sex": "femal |
2153
            2154
2155
        # Attempt login with email as username.
2156
2157
        file_system_username =
2158
             '[susanmiller@gmail.com](mailto:susanmiller@gmail.com)"
2159
```

```
2160
        file_system_password = [a["password"] for a in
2161
        → apis.supervisor.show_account_passwords()
2162
        if a["account_name"] == "file_system"][0]
2163
        file_system_login_response = apis.file_system.login(
2164
        username=file_system_username, password=file_system_password)
        print(file_system_login_response)
2165
2166
2167
           { "access_token": "eyJhbGciOiJIUzI1NiIsInR5cCI6IkpXVCJ9... < REDACTE |
           → D>...", "token_type": "Bearer"}
2168
2169
2170
        # Try to set token manually (guessing an API), then list directory.
2171
2172
        apis.file_system.set_access_token("...<REDACTED>...")
2173
        downloads_path = '~/downloads'
2174
        downloads_files = apis.file_system.show_directory(path=downloads_path)
2175
        print (downloads_files)
2176
2177
           Execution failed. Traceback:
2178
           Exception: No API named 'set_access_token' found in the
2179
           \hookrightarrow file_system app.
2180
2181
2182
        # Call show_directory again, assuming session is persisted after
2183
        → login.
2184
        downloads_path = '~/downloads'
2185
        downloads_files = apis.file_system.show_directory(path=downloads_path)
2186
        print (downloads_files)
2187
2188
           Execution failed. Traceback:
2189
2190
           Exception: Response status code is 401:
2191
           {"message":"You are either not authorized to access this

→ file_system API endpoint or your access token is missing,

2192
           → invalid or expired."}
2193
2194
2195
        # Retry: fresh login, try immediate call.
2196
2197
        login_response = apis.file_system.login(
2198
        username=file_system_username, password=file_system_password)
2199
        access_token = login_response.get("access_token")
        print(access_token) # for debugging
2200
2201
           eyJhbGciOiJIUzI1NiIsInR5cCI6IkpXVCJ9...<REDACTED>...
2202
2203
2204
        # Try show_directory once more after login.
2205
2206
        downloads_path = '~/downloads'
2207
        downloads_files = apis.file_system.show_directory(path=downloads_path)
2208
        print (downloads_files)
2209
2210
```

```
2214
2215
           Execution failed. Traceback:
2216
2217
           Exception: Response status code is 401:
2218
           {"message":"You are either not authorized to access this

→ file_system API endpoint or your access token is missing,

2219
           → invalid or expired."}
2220
2221
2222
        # Inspect file_system APIs again for any token/session helpers.
2223
2224
        print(apis.api_docs.show_api_descriptions(app_name='file_system'))
2225
2226
           [{"name": "show_account", "description": "Show your account
2227
              information. Unlike show_profile, this includes private
2228
               information."},
           {"name": "signup", "description": "Sign up to create account."},
2229
           {"name": "delete_account", "description": "Delete your account."},
2230
           {"name": "update_account_name", "description": "Update your first or
2231
           → last name in the account profile."},
2232
           {"name":"login", "description": "Login to your account."},
2233
           . . .
           {"name": "copy_directory", "description": "Copy a directory to
2234

    another location."
},
2235
           {"name": "move_directory", "description": "Move a directory to
2236

    another location."},
2237
           {"name": "compress_directory", "description": "Compress a
2238

    directory."},
           {"name": "decompress_file", "description": "Decompress a compressed
2239

  file."
}]
2240
2241
2242
        # Try logout -> login -> call again.
2243
2244
        apis.file_system.logout()
2245
        login_response = apis.file_system.login(
2246
        username=file_system_username, password=file_system_password)
        downloads_path = '~/downloads'
2247
        downloads_files = apis.file_system.show_directory(path=downloads_path)
2248
        print (downloads_files)
2249
2250
           Execution failed. Traceback:
2251
2252
           Exception: Response status code is 422:
2253
           {"message": "You are either not authorized to access this

→ file_system API endpoint or your access token is missing,

2254
           \hookrightarrow invalid or expired."}
2255
2256
2257
        # Final retry without logout; still unauthorized.
2258
2259
        login_response = apis.file_system.login(
2260
        username=file_system_username, password=file_system_password)
2261
        downloads_path = '~/downloads'
        downloads_files = apis.file_system.show_directory(path=downloads_path)
2262
        print (downloads_files)
2263
2264
```

```
2268
2269
          Execution failed. Traceback:
2270
2271
          Exception: Response status code is 401:
2272
           {"message":"You are either not authorized to access this

→ file_system API endpoint or your access token is missing,

2273
           → invalid or expired."}
2274
2275
2276
        # Supervisor cannot manage files; mark task failed due to
2277
        \hookrightarrow authorization issues.
2278
2279
        print(apis.api_docs.show_api_descriptions(app_name='supervisor'))
2280
2281
2282
           { "name": "show_active_task", "description": "Show the currently
           → active task assigned to you by the supervisor."},
2283
           {"name": "complete task", "description": "Mark the currently active
2284

→ task as complete with the given answer."},
2285
           {"name": "show_profile", "description": "Show your supervisor's
2286
           → profile information."},
           {"name": "show_addresses", "description": "Show your supervisor's
2287

    addresses."
},
2288
           { "name": "show_payment_cards", "description": "Show your
2289
           → supervisor's payment_cards."},
2290
           { "name": "show_account_passwords", "description": "Show your
2291
           ]
2293
2294
        apis.supervisor.complete_task(status="fail")
2295
```

Example E.3: AppWorld trajectory example - gpt-4.1-mini agent successes with history compression

Problem: Delete all .pdf files from ~/downloads.

```
2322
2323
           (... previous raw history is omitted for brevity ...)
2324
           <history_summary>
2325
2326
           1. REASONING
               - User requested deletion of all .pdf files in ~/downloads.
2327
              - file_system APIs require authentication; initial calls
2328

→ failed with 401 Unauthorized.

2329
              (... omitted ...)
2330
              - Confirmed login API returns access_token but no mechanism to
2331
              \hookrightarrow use it in subsequent calls found.
              - Need to find a way to authorize file_system API calls with
2332
               \rightarrow the access token or alternative approach.
2333
2334
           2. VARS
2335
              | name
                                        | value
2336
               \hookrightarrow | purpose
              |----|---
                              -----|
2337
              | file_system_username | "susanmiller@gmail.com"
2338

→ | Username for file_system login |
2339
              | file_system_password | "8nNw!jZ"
2340
               \hookrightarrow | Password for file_system login
2341
              | file_system_token
               → "eyJhbGciOiJIUzI1NiIsInR5cCI6IkpXVCJ9...sBA8nTOB..." |
2342
               \hookrightarrow Access token from login |
2343
               | downloads_path | "~/downloads"
2344

→ | Target directory for file listing |
2345
           3. TODO
2346
              - Determine how to pass or set the access token for
2347
              \rightarrow file_system API calls to authorize them.
2348
              - Once authorized, list files in ~/downloads to identify .pdf
2349
               \hookrightarrow files.
2350
              - Delete all .pdf files found in ~/downloads.
2351
           4. COMPLETED
2352
              - Retrieved file_system account password from supervisor app.
2353
              - Successfully logged in to file_system app with username and
2354

→ password, obtained access token.

2355
              - Confirmed file_system APIs require authorization; calls
2356
               \hookrightarrow without token fail.
              - Verified no supervisor app API to proxy file system
2357
               → operations.
2358
              - Confirmed login API returns access token but no visible
2359
               \rightarrow method to use it in subsequent calls.
2360
2361
           5. GUARDRAILS
              - Always provide both username and password for login API.
2362
              - Access tokens must be explicitly handled or passed if
2363

→ environment does not auto-use them.

2364
              - Validate authorization before calling protected APIs to
2365

→ avoid 401 errors.

2366
              - Avoid redundant login attempts without token usage

→ resolution.

2367
              - Preserve tokens and credentials securely for session
2368

→ continuity.

2369
              </HISTORY_SUMMARY>
2370
2371
2372
```

The show_directory API requires the access_token as a parameter.

2373

```
2376
        # List files in "~/downloads" using the token.
2377
2378
        downloads_path = "~/downloads"
2379
        file_system_access_token =
2380
        → "eyJhbGciOiJIUzI1NiIsInR5cCI6IkpXVCJ9...sBA8nTOB..."
        print(apis.file_system.show_directory(
2382
        access_token=file_system_access_token,
2383
        directory_path=downloads_path,
2384
        entry_type="files",
2385
        recursive=False))
2386
2387
2388
           "/home/susan/downloads/aerobics_class.jpg",
           "/home/susan/downloads/art_inspiration_sketches.zip",
2389
           "/home/susan/downloads/bakeware.pdf",
2390
           "/home/susan/downloads/bowls_for_snacks.pdf",
2391
           "/home/susan/downloads/candy.pdf",
2392
           "/home/susan/downloads/charcoal_pencils.pdf",
2393
           . . . ,
           "/home/susan/downloads/travel_insurance.pdf",
2394
           "/home/susan/downloads/workout_clothes.pdf",
2395
           "/home/susan/downloads/world_cuisine_recipe_book.epub",
2396
           "/home/susan/downloads/world_landmarks_photo_album.zip"
2397
2398
2399
2400
        # Filter to only ".pdf" files.
2401
        downloads_files = [
2402
        "/home/susan/downloads/aerobics_class.jpg",
2403
        "/home/susan/downloads/bakeware.pdf",
2404
        "/home/susan/downloads/bowls_for_snacks.pdf",
2405
        "/home/susan/downloads/candy.pdf",
        "/home/susan/downloads/charcoal_pencils.pdf",
2406
2407
        "/home/susan/downloads/travel_insurance.pdf",
2408
        "/home/susan/downloads/workout_clothes.pdf"
2409
2410
        pdf_files = [p for p in downloads_files if p.lower().endswith(".pdf")]
        print(pdf_files)
2411
2412
2413
           "/home/susan/downloads/bakeware.pdf",
2414
           "/home/susan/downloads/bowls_for_snacks.pdf",
2415
           "/home/susan/downloads/candy.pdf",
2416
           "/home/susan/downloads/charcoal_pencils.pdf",
2417
           ...,
2418
            /home/susan/downloads/travel insurance.pdf",
           "/home/susan/downloads/workout_clothes.pdf"
2419
2420
2421
2422
        # Confirm delete_file API parameters before deletion.
2423
2424
        print (apis.api_docs.show_api_doc(app_name="file_system",
2425

    api_name="delete_file"))

2426
2427
```

```
2430
2431
2432
            "app_name": "file_system",
2433
            "api_name": "delete_file",
            "path":"/file",
2434
            "method": "DELETE",
2435
           "description": "Delete a file.",
2436
            "parameters":
2437
            {"name": "file_path", "type": "string", "required":true,
2438
           "description": "Path of the file. Path can be absolute, starting
2439
            \hookrightarrow with '/', or relative to the user's home directory, starting
            ⇔ with '~/'."},
2440
            {"name":"access_token","type":"string","required":true,
2441
           "description": "Access token obtained from file_system app
2442
            → login."}
2443
           "response_schemas":{
2444
           "success": { "message": "string", "file_path": "string" },
2445
           "failure": { "message": "string" }
2446
2447
           }
2448
2449
2450
        # Delete each PDF using the access_token.
2451
        delete_results = []
2452
        for pdf in pdf_files:
2453
        res = apis.file_system.delete_file(
2454
        file_path=pdf, access_token=file_system_access_token)
2455
        delete_results.append(res)
2456
        # Mark task complete.
2457
2458
        apis.supervisor.complete_task()
2459
2460
```

```
2484
       Compression E.1: AppWorld History Compression Example
2485
2486
        task: I got some Splitwise group invitations over phone text
2487
           messages yesterday.
2488
        If their number is in my phone contact book, accept it, otherwise
2489
           delete those messages.
2490
2491
2492
        History Compression (Prompting baseline):
2493
        ### REASONING
2494
       The agent began by identifying relevant apps for handling text
2495
           messages and Splitwise invitations, then explored the phone app
2496
           's APIs to manage messages and contacts. To search for
2497
           Splitwise group invitations received via text message, the
2498
           agent determined it needed to log in to the phone app to obtain
2499
            an access token. Using credentials from the supervisor app,
           the agent successfully logged in and retrieved the access token
2500
           . The agent then calculated yesterday's date range and searched
2501
            for text messages containing "Splitwise," retrieving a list of
2502
            messages. The rationale for retaining message details (sender
2503
           name, phone number, contact_id, message content, sent_at) is to
2504
            enable the next step: checking if the sender's number is in
           the contact book, so the agent can accept or delete invitations
2505
            as instructed. Earlier steps (API exploration, login, date
2506
           calculation) directly enabled the message search and data
2507
           extraction.
2509
        ### COMPLETED
2510
        - Identified relevant apps and APIs for the task.
2511
        - Successfully logged in to the phone app and obtained access token
2512
2513
        - Calculated yesterday's date range.
2514
        - Searched and retrieved text messages mentioning "Splitwise" from
           yesterday and beyond, including sender details and message
2515
           content.
2516
2517
2518
2519
       History Compression (ACON (utility step)):
       1. REASONING
           - Identified "phone" app for managing text messages and contacts
2521
2522
           - Retrieved API docs to find endpoints for searching messages
2523
           and contacts.
2524
           - Determined yesterday's date (May 17, 2023) for filtering
           relevant messages.
2525
           - Checked "search_text_messages" API spec for required
2526
           parameters.
2527
           - Used supervisor app to obtain phone app password.
2528
           - Logged in to phone app to get access_token.
2529
           - Searched for messages containing "Splitwise" with page_limit
           =20, sorted by -created_at.
2530
           - Located Splitwise group invitation messages; sender contact_id
2531
            used to determine if sender is in contact book.
2532
2533
       2. VARS
2534
                          | value
          | name
                                   | purpose
2535
2536
           2537
```

```
2538
           | access_token | eyJhbGciOiJIUzI1NiIsInR5cCI6IkpXVCJ9.
2539
           eyJzdWIiOiJwaG9uZSs3OTM0ODE4MTkxIiwiZXhwIjoxNjg0NDEyMDk4fQ.
2540
            ad3t3S_Xj7YhoDDFZeW4BlVL4dNniMdfaXC1143fbzw | Auth for phone
2541
            app APIs
2542
           | phone_number | 7934818191
                                    | User's phone number
2543
2544
           | phone_password| X53y6aW
2545
                                   | Phone app password
2546
           | search_query | "Splitwise"
2547
                                    | Filter for relevant messages
2548
2549
                            | -created_at
           | sort_by
2550
                                    | Get latest messages first
2551
2552
           | page_index
                            | 0
                                    | Current page index for pagination
2553
2554
           | page_limit
                           1 20
2555
                                    | Number of results per page
2556
2557
                           | 2023-05-17
           | target_date
                                    | Date to filter messages
2558
2559
           | invitation_msgs | [IDs: 16794, 16800]
2560
                                   | Splitwise invitations from 2023-05-17
2561
           | sender_info | [{id:1857, name: "Alex White", phone
            :"3707865891"}, {id:null, name:"Anita Burch", phone
2563
            :"3643463570"}] | For contact book check
2564
2565
        3. TODO
2566
           - For each Splitwise invitation message from 2023-05-17:
2567
             - Check if sender's phone number exists in contact book (
2568
           search contacts API).
             - If in contact book, accept invitation (Splitwise API, if
2569
           needed).
2570
             - If not in contact book, delete the message (
2571
            delete_text_message API).
2572
           - Paginate search_text_messages if more relevant messages may
           exist.
2573
2574
        4. COMPLETED
2575
           - Identified relevant apps and APIs for task.
2576
           - Retrieved phone app password from supervisor.
2577
           - Logged in to phone app and obtained access_token.
           - Searched and filtered text messages for Splitwise invitations
2578
           on target date.
2579
           - Isolated relevant message IDs and sender info for next steps.
2580
2581
        5. GUARDRAILS
2582
           - Always re-declare VARS after memory reset.
           - Paginate message search until no more results.
2583
           - Validate API parameters before each call.
2584
           - Avoid redundant logins and API doc look-ups.
2585
           - Only act on messages matching both query and target date.
2586
2587
2588
2589
        History Compression (ACON (utility step + compression step)):
2590
2591
```

2592 1. REASONING 2593 I identified voice messages from the day before yesterday (2023-05-16) containing Splitwise group invitations. To follow 2595 instructions, I need to check if the senders' phone numbers are 2596 in Sierra's contact book before accepting or deleting these messages. 2598 2. VARS access_token=<token>, splitwise_invite_voice_ids 2600 =[17022,17036,17028,17034], sender_numbers 2601 =['9595488531','7375474944','2123782445','3808719492'] 2602 3. TODO 2603 - Search Sierra's contact book for sender_numbers. 2604 - Accept Splitwise invitations from contacts found. 2605 - Delete voice messages from non-contacts. 2606 4. COMPLETED 2607 - Logged in to phone app and obtained access_token. 2608 - Retrieved voice messages from 2023-05-16. 2609 - Identified Splitwise invitation messages and sender numbers. 2610 2611 5. GUARDRAILS Only process Splitwise invitations from contacts; delete others. 2612 2613 2614 Compression E.2: 8-objective QA History Compression Example 2615 2616 task: when was the last time the dodgers played yankees in the 2617 world series?; where will be the next olympics be held?; who is 2618 tinker air force base named after?; who does the voice of cruz 2619 on cars 3?; who's the original singer of help me make it through the night?; where does the movie mothers day take place 2620 ?; what is the big gold dome in jerusalem?; when did russia 2621 join the world economic forum? 2622 2623 2624 History Compression (Prompting baseline): 2625 ### REASONING 2626 The agent began by addressing the user's multi-part query, starting 2627 with "when was the last time the Dodgers played Yankees in the 2628 World Series?" Through Wikipedia searches, the agent 2629 identified 1981 as the most recent World Series meeting between the Dodgers and Yankees. This answer was retained due to clear 2630 evidence in the retrieved documents. The agent then attempted 2631 to answer "where will be the next Olympics be held?" but 2632 initial searches returned irrelevant or outdated results (e.g., 2633 Amsterdam 1928, Buenos Aires 2018). Recognizing the lack of 2634 direct answers, the agent refined its search queries multiple times, but still did not retrieve the current location for the 2635 next Olympics. The agent's iterative approach demonstrates a 2636 commitment to accuracy and highlights the need to adjust search 2637 strategies when initial attempts fail. Retaining details of 2638 failed searches helps prevent repeated mistakes and guides 2639 future queries. 2640 ### COMPLETED 2641

49

- Identified the last time the Dodgers played the Yankees in the

World Series: 1981.

2642

```
History Compression (ACON (utility step)):
2647
2648
        ### REASONING
2649
        - Used Wikipedia search for each question; refined Olympics query
2650
           after initial results lacked direct info.
        - Collapsed repeated/irrelevant Olympic bid results; focused on
2651
           confirmed host city.
2652
        - For Tinker AFB, confirmed name origin from first document.
2653
2654
        ### COMPLETED
2655
        - Last Dodgers vs Yankees World Series: 1955 (source: "1955 World
2656
            Series")
        - Next Olympics location: Paris, 2024 (source: "Budapest bid for
2657
           the 2024 Summer Olympics")
2658
        - Tinker Air Force Base named after Major General Clarence L.
2659
           Tinker (source: "Tinker Air Force Base")
2660
        ### STATE LOG
2661
       FACT CACHE:
2662
        - When was the last time the Dodgers played Yankees in the World
2663
           Series? -> 1955 ("1955 World Series")
2664
        - Where will the next Olympics be held? -> Paris, 2024 ("Budapest
2665
           bid for the 2024 Summer Olympics")
        - Who is Tinker Air Force Base named after? -> Major General
2666
           Clarence L. Tinker ("Tinker Air Force Base")
2667
2668
       PENDING:
2669
        - Who does the voice of Cruz on Cars 3?
2670
        - Who is the original singer of "Help Me Make It Through the Night
            "?
2671
        - Where does the movie "Mother's Day" take place?
2672
        - What is the big gold dome in Jerusalem?
2673
        - When did Russia join the World Economic Forum?
2674
2675
       ATTEMPTED QUERIES / LOOP GUARD:
        - "Dodgers Yankees World Series history" OK
2676
        - "next Olympics location" FAIL (no direct info)
2677
        - "2024 Olympics location" OK
        - "Tinker Air Force Base name origin" OK
2679
        - "next Olympics location" tried twice (loop guard: must
2680
           reformulate/abandon)
2681
        ERROR NOTES:
2682
        None
2683
2684
        (End of summary - no extra text)
2685
2686
2687
        History Compression (ACON (utility step + compression step)):
2688
2689
        ### FACT CACHE
2690
        - When was the last time the Dodgers played Yankees in the World
            Series? -> 1956
2691
2692
        ### PENDING
2693
        where will the next olympics be held?, who is tinker air force base
2694
            named after?, who does the voice of cruz on cars 3?, who is
2695
            the original singer of help me make it through the night?,
           where does the movie mothers day take place?, what is the big
2696
            gold dome in jerusalem?, when did russia join the world
2697
            economic forum?
2698
2699
```

```
2700
        ### BLOCKED
2701
        next Olympics location: 3x failed searches (loop-quard)
2702
2703
        ### REASONING
2704
        Olympics location queries returned irrelevant results; alternate
            sources or query phrasing needed.
2705
2706
2707
        Compression E.3: OfficeBench History Compression Example
2708
2709
        task: read company revenues, save the top 2 revenue years in top2.
2710
            xlsx, generate reports.docx
2711
2712
2713
        History Compression (Prompting baseline):
2714
        ### REASONING
2715
2716
        The agent began by identifying the relevant data file (
2717
            company_revenues.xlsx) containing company revenue information.
2718
           After reading the file, the agent determined the top 2 revenue
           years: 2014 ($4,983,791) and 2015 ($4,611,210). The agent then
2719
            created a new Excel file (top2.xlsx) and started populating it
2720
            with the required data, first adding headers ("Year", "Revenue
2721
            ($)") and then entering the first top year (2014) in the second
2722
            row, first column. The rationale for retaining row and column
2723
            indices, file paths, and the top 2 years is to ensure accurate
            and complete transfer of data and to facilitate the next steps
2724
            (completing the top2.xlsx file and generating reports.docx).
2725
2726
        ### COMPLETED
2727
        - Located and read company_revenues.xlsx.
2728
        - Identified top 2 revenue years: 2014 ($4,983,791), 2015 ($4
2729
            ,611,210).
2730
        - Created top2.xlsx and added headers ("Year", "Revenue ($)").
2731
        - Entered "2014" in (2, 1) of top2.xlsx.
2732
2733
2734
        History Compression (ACON (utility step)):
2735
2736
        ### REASONING
2737
        - Located company revenue data in company_revenues.xlsx.
        - Read all rows; identified top 2 revenue years: 2014 ($4,983,791)
2738
            and 2015 ($4,611,210).
2739
        - Created top2.xlsx and began writing header and first top year.
2740
        - Each numeric value was copied verbatim from source.
2741
        - No errors or mismatches so far; headers and first year written
2742
            successfully.
        - Pending: finish writing revenue for 2014, add 2015 row, then
2743
            generate reports.docx.
2744
2745
        Verify numeric values directly from source file before use.
2746
2747
2748
        ### DATA STATE
2749
2750
        #### company_revenues.xlsx (source table)
2751
        | Row | Year | Revenue ($) |
2752
        |-----|
        | 2 | 2004 | 4439044
2753
```

```
2754
        | 3
             | 2005 | 2212203
2755
       | 4
             | 2006 | 1193225
2756
       | 5
             | 2007 | 2793265
2757
       | 6
            | 2008 | 2080460
2758
            | 2009 | 1791112
       | 7
       | 8
             | 2010 | 3098049
2759
       1 9
             | 2011 | 3975966
2760
       | 10
             | 2012 | 3251777
2761
             | 2013 | 2918667
       | 11
2762
       | 12
             | 2014 | 4983791
2763
             | 2015 | 4611210
       | 13
       | 14
             | 2016 | 3294631
2764
       | 15
             | 2017 | 1463500
2765
             | 2018 | 3945846
       | 16
2766
       | 17
             | 2019 | 3907271
2767
       | 18
             | 2020 |
                       3651262
2768
       | 19
             | 2021 |
                       2005820
       | 20
             | 2022 | 2170828
2769
       | 21 | 2023 | 2832545
2770
2771
       #### top2.xlsx (current state)
2772
       | (row, col) | Value
2773
       |-----
       | (1, 1)
                   | Year
2774
       | (1, 2)
                    | Revenue ($)
2775
       (2, 1)
                    | 2014
2776
       (2, 2)
                    | [Empty Cell]
2777
2778
        ___
2779
       ### TODO / PENDING
2780
       - Write 2014 revenue (4983791) to (2,2) in top2.xlsx.
2781
        - Write 2015 year (3,1) and revenue (4611210) (3,2) in top2.xlsx.
2782
         Generate reports.docx summarizing top 2 revenue years.
2783
       - Re-read source rows before writing each numeric value for
2784
           consistency.
2785
2786
2787
       ### COMPLETED
2788
       - Located and read company_revenues.xlsx.
       - Created top2.xlsx.
2789
       - Wrote headers and 2014 year to top2.xlsx.
2790
2791
```

Compression E.4: AppWorld Observation Compression Example

2792

27932794

2795

2796

2797

2798

279928002801

2802 2803 2804

2805

2806

```
task: I am going for a 15-minute drive without internet. Play an
    album from my Spotify library that already has enough
    downloaded songs for it, so I do not have to repeat.

action: print(apis.api_docs.show_api_descriptions(app_name='spotify
    '))

Raw Observation:
[
{
    "name": "show_account",
    "description": "Show your account information. Unlike
    show_profile, this includes private information."
},
```

```
2808
2809
          "name": "signup",
2810
          "description": "Sign up to create account."
2811
2812
          "name": "delete_account",
2813
          "description": "Delete your account."
2814
2815
2816
          "name": "update_account_name",
2817
          "description": "Update your first or last name in the account
            profile."
2818
2819
2820
          "name": "login",
2821
          "description": "Login to your account."
2822
         },
2823
          "name": "logout",
2824
          "description": "Logout from your account."
2825
2826
2827
          "name": "send_verification_code",
          "description": "Send account verification code to your email
2828
            address."
2829
         },
2830
2831
          "name": "verify_account",
          "description": "Verify your account using the verification code
            sent to your email address."
2833
2834
2835
          "name": "send_password_reset_code",
2836
          "description": "Send password reset code to your email address."
2837
2838
          "name": "reset_password",
2839
          "description": "Reset your password using the password reset code
2840
             sent to your email address."
2841
2842
          "name": "show_profile",
2843
          "description": "Show public profile information of a user."
2844
2845
2846
          "name": "show_genres",
2847
          "description": "Show the list of all music genres."
2848
2849
          "name": "search_songs",
2850
          "description": "Search for songs with a query."
2851
2852
          "name": "show_song",
2853
          "description": "Get details of a specific song."
2854
2855
2856
          "name": "show_song_privates",
2857
          "description": "Show information about the song that is private
            to the user."
2858
         },
2859
2860
          "name": "like_song",
2861
```

```
"description": "Like a song."
2863
         },
2864
2865
          "name": "unlike_song",
2866
          "description": "Unlike a song."
2868
          "name": "show_liked_songs",
          "description": "Get a list of songs you have liked."
2870
2871
          "name": "search_albums",
2872
          "description": "Search for albums with a query."
2873
2874
2875
          "name": "show_album",
2876
          "description": "Get details of a specific album."
2877
2878
          "name": "show_album_privates",
2879
          "description": "Show information about the album that is private
2880
            to the user."
2881
2882
          "name": "like_album",
2883
          "description": "Like a album."
2884
2885
          "name": "unlike_album",
          "description": "Unlike a album."
2887
2888
2889
          "name": "show_liked_albums",
2890
          "description": "Get a list of albums you have liked."
2891
2892
          "name": "show_playlist_library",
2893
          "description": "Get a list of playlists in the user's playlist
2894
            library."
2895
2896
          "name": "search_playlists",
2897
          "description": "Search for playlists with a query. It will search
             over all public playlists and your own private playlists."
2899
         },
2900
2901
          "name": "create_playlist",
          "description": "Create a new playlist."
2902
2903
2904
          "name": "show_playlist",
2905
          "description": "Get detailed information about a specific
2906
            playlist. You can view your own playlists or others' playlists
            if they are public."
2907
         },
2908
2909
          "name": "delete_playlist",
2910
          "description": "Delete a playlist."
2911
         },
2912
          "name": "update_playlist",
2913
          "description": "Update a playlist title or privacy."
2914
2915
```

```
2916
2917
          "name": "show_playlist_privates",
2918
          "description": "Show information about the playlist that is
2919
            private to the user."
2920
2921
          "name": "like_playlist",
2922
          "description": "Like a playlist."
2923
2924
2925
          "name": "unlike_playlist",
          "description": "Unlike a playlist."
2926
2927
2928
          "name": "show_liked_playlists",
2929
          "description": "Get a list of playlists you have liked."
2930
         },
2931
          "name": "search_artists",
2932
          "description": "Search for artists with a query."
2933
2934
2935
          "name": "show_artist",
          "description": "Get details of a specific artist."
2936
         }.
2937
2938
          "name": "show_artist_following",
2939
          "description": "Show if the user is following the artist."
2940
2941
          "name": "show_song_library",
2942
          "description": "Get a list of songs in the user's song library."
2943
2944
2945
          "name": "add_song_to_library",
          "description": "Add a song to the user's song library."
2946
2947
2948
          "name": "remove_song_from_library",
2949
          "description": "Remove a song from the user's song library."
2950
         },
2951
          "name": "show_album_library",
2952
          "description": "Get a list of albums in the user's album library
2953
            . "
2954
2955
          "name": "add_album_to_library",
2956
          "description": "Add an album to the user's album library."
2957
2958
2959
          "name": "remove_album_from_library",
2960
          "description": "Remove an album from the user's album library."
2961
2962
          "name": "add_song_to_playlist",
2963
          "description": "Add a song to a playlist."
2964
2965
          "name": "remove_song_from_playlist",
2966
          "description": "Remove a song from a playlist."
2967
2968
2969
```

```
2970
          "name": "show_downloaded_songs",
2971
          "description": "Get a list of downloaded songs."
2972
         },
2973
2974
          "name": "download_song",
          "description": "Download a song."
2975
2976
2977
          "name": "remove_downloaded_song",
2978
          "description": "Remove a song from downloads."
2979
         },
2980
          "name": "show_following_artists",
2981
          "description": "Get a list of artists the user is following."
2982
         },
2983
2984
          "name": "follow_artist",
          "description": "Follow an artist."
2985
2986
2987
          "name": "unfollow_artist",
2988
          "description": "Unfollow an artist."
2989
2990
          "name": "show_song_reviews",
2991
          "description": "Get a list of reviews for a song."
2992
2993
2994
          "name": "review_song",
          "description": "Rate or review a song."
2995
2996
2997
          "name": "show_song_review",
2998
          "description": "Show a song review."
2999
3000
          "name": "delete_song_review",
3001
          "description": "Delete a song review."
3002
3003
3004
          "name": "update_song_review",
          "description": "Update a song review."
3005
         },
3006
3007
          "name": "show_album_reviews",
3008
          "description": "Get a list of reviews for an album."
3009
         },
3010
          "name": "review_album",
3011
          "description": "Rate or review an album."
3012
         },
3013
3014
          "name": "show_album_review",
          "description": "Show an album review."
3015
         },
3016
3017
          "name": "delete_album_review",
3018
          "description": "Delete an album review."
3019
         },
3020
          "name": "update_album_review",
3021
          "description": "Update an album review."
3022
3023
```

```
3024
3025
          "name": "show_playlist_reviews",
          "description": "Show a list of reviews for your playlist or
3027
            others' public playlist."
3028
3029
          "name": "review_playlist",
3030
          "description": "Rate or review a playlist."
3031
3032
3033
          "name": "show_playlist_review",
          "description": "Show a playlist review."
3034
3035
3036
          "name": "delete_playlist_review",
3037
          "description": "Delete a playlist review."
3038
         },
3039
          "name": "update_playlist_review",
3040
          "description": "Update a playlist review."
3041
         },
3042
3043
          "name": "show_payment_cards",
          "description": "Get a list of users payment cards."
3044
         }.
3045
3046
          "name": "add_payment_card",
3047
          "description": "Add a new payment card."
3048
         },
3049
          "name": "show_payment_card",
3050
          "description": "Get details of a payment card."
3051
3052
3053
          "name": "delete_payment_card",
          "description": "Delete payment card information."
3054
3055
3056
          "name": "update_payment_card",
3057
          "description": "Update payment card information."
3058
3059
          "name": "show_current_song",
3060
          "description": "Show details of the current song on the queue."
3061
         },
3062
3063
          "name": "play_music",
          "description": "Play music based on various criteria. You can
3064
            pass, at most, any one of queue_position, song_id, album_id or
3065
            playlist_id. If one of song_id, album_id or playlist_id is
3066
            passed, that song, album or playlist will be added to the queue
3067
            and played. Otherwise, the queue will remain unchanged. If
3068
            queue_position is passed, the song at that position in the
            queue will be played. If none is passed, the current song in
3069
            the queue will be played."
3070
         },
3071
3072
          "name": "pause_music",
3073
          "description": "Pause the currently playing song."
3074
3075
          "name": "previous_song",
3076
          "description": "Go to the previous song in the song queue."
3077
```

```
3078
         },
3079
3080
          "name": "next_song",
3081
          "description": "Go to the next song in the song queue."
3082
3083
          "name": "move_song_in_queue",
3084
          "description": "Move a song in the queue to a new position."
3085
3086
3087
          "name": "seek_song",
          "description": "Seek the current song to the given number of
3088
            seconds."
3089
3090
3091
          "name": "loop_song",
3092
          "description": "Set whether to loop the current song."
3093
3094
          "name": "shuffle song queue",
3095
          "description": "Shuffle songs in the music player queue."
3096
3097
          "name": "show_song_queue",
3098
          "description": "Get the music player song queue. Songs are played
3099
             in the order of the queue in a cycle."
3100
3101
3102
          "name": "add_to_queue",
          "description": "Add a song, album or playlist to the music player
3103
             song queue."
3104
         },
3105
3106
          "name": "clear_song_queue",
3107
          "description": "Clear the music player song queue."
3108
3109
          "name": "remove_song_from_queue",
3110
          "description": "Remove a song at the given position from the
3111
            music player song queue."
3112
         },
3113
          "name": "show_volume",
3114
          "description": "Get the volume level of the music player."
3115
         },
3116
3117
          "name": "set_volume",
          "description": "Set the volume level of the music player."
3118
3119
3120
          "name": "show_recommendations",
3121
          "description": "Get personalized song recommendations for the
3122
            user."
3123
         },
3124
          "name": "show_premium_plans",
3125
          "description": "Show information about premium plans available."
3126
3127
          "name": "show_premium_subscriptions",
3128
          "description": "Show your premium subscription history."
3129
3130
3131
```

```
3132
          "name": "subscribe_premium",
3133
          "description": "Subscribe to premium membership."
3134
3135
3136
          "name": "download_premium_subscription_receipt",
          "description": "Download the receipt for a premium subscription."
3137
3138
3139
3140
3141
3142
        Obs. Compression (Prompting baseline):
3143
        The Spotify API provides:
3144
        - show_album_library: get user's album library.
3145
        - show_downloaded_songs: get list of downloaded songs.
3146
        - show_album: get details of a specific album.
3147
3148
3149
        Obs. Compression (ACON (utility step)):
3150
3151
            "name": "show_album_library",
3152
            "description": "Get a list of albums in the user's album
3153
            library."
3154
          },
3155
          {
3156
            "name": "show_downloaded_songs",
            "description": "Get a list of downloaded songs."
3157
          },
3158
3159
            "name": "show_album",
3160
            "description": "Get details of a specific album."
3161
3162
            "name": "play_music",
3163
            "description": "Play music based on various criteria. You can
3164
            pass, at most, any one of queue_position, song_id, album_id or
3165
            playlist_id. If one of song_id, album_id or playlist_id is
3166
            passed, that song, album or playlist will be added to the queue
             and played. Otherwise, the queue will remain unchanged. If
3167
            queue_position is passed, the song at that position in the
3168
            queue will be played. If none is passed, the current song in
3169
            the queue will be played."
3170
3171
        1
3172
3173
3174
3175
        History Compression (ACON (utility step + compression step)):
3176
        [{"name": "show_album_library", "description": "Get user's album
3177
            library."}, { "name": "show_downloaded_songs", "description": "Get
3178
            downloaded songs."}, { "name": "show_album_privates", "description
3179
            ": "Show album private info." }, { "name": "play_music", "description
3180
            ":"Play music; album_id allowed."}]
3181
3182
```