Mobile-Agent-E: Self-Evolving Mobile Assistant for Complex Tasks

Zhenhailong Wang¹, Haiyang Xu², Junyang Wang², Xi Zhang², Ming Yan², Ji Zhang², Fei Huang², Heng Ji¹

¹University of Illinois Urbana-Champaign ²Alibaba Group {wangz3,hengji}@illinois.edu, shuofeng.xhy@alibaba-inc.com



Figure 1: Introducing Mobile-Agent-E, a novel hierarchical agentic framework designed for complex real-world mobile tasks. Mobile-Agent-E disentangles high-level planning from low-level actions, significantly outperforming previous state-of-the-art approaches [38, 29, 28]. Equipped with a novel self-evolution module that learns general *Tips* and reusable *Shortcuts* from past experiences, Mobile-Agent-E further enhances both performance and efficiency.

Abstract

Recent advancements in large multimodal model (LMM)-based mobile agents have demonstrated promising capabilities in acting within mobile environments. However, current approaches face significant limitations: (1) they fall short in addressing real-world human needs, which involve complex, open-ended, and reasoning-intensive tasks; and (2) they lack mechanisms to learn and improve from prior experiences. To address these challenges, we introduce Mobile-Agent-E, a hierarchical agentic framework capable of self-evolution through past experience. Mobile-Agent-E adopts a multi-level communication protocol for reasoning, perception, and error recovery, explicitly separating high-level planning from low-level action decisions. It also introduces a novel self-evolution module that maintains a persistent long-term memory comprising Tips and Shortcuts, enabling continual refinement of task performance and efficiency. To bridge the gap in existing benchmarks for complex, open-ended tasks, we further present a new benchmark, Mobile-Eval-E, alongside a new evaluation metric, the Satisfaction Score. Empirical results show that Mobile-Agent-E achieves a 22% absolute improvement over previous state-of-the-art approaches across three LMM backbones. We also

provide a comprehensive analysis of the impact of the self-evolution mechanism and outline promising directions for future work.

1 Introduction

Recent advancements in large multimodal models (LMMs) [18, 2, 26] have led to the emergence of LMM-based GUI agents [30, 17] capable of acting in the Web, PC, and mobile environments. Despite these initial successes, current research on mobile agents [29, 38, 28, 11] has yet to fully address the challenges of real-world mobile tasks. We identify two key limitations below.

First, existing mobile agents and benchmarks focus primarily on *goal-oriented* tasks, such as "Create a new contact for {name}. Their number is {number}" [20]. These tasks typically follow a linear ground-truth trajectory and have a single success state. However, we argue that tasks more representative of real human needs are significantly more complex. They often require: (1) intensive reasoning to satisfy multiple constraints; (2) long-horizon planning that spans across multiple apps; and (3) open-ended exploration, where vague instructions demand active information gathering. For instance, as illustrated in Figure 1, online shopping often involves navigating across different apps to compare prices and find the best deal.

Second, unlike humans who quickly adapt to recurring tasks on new devices, current mobile agents lack the ability to learn from prior experiences. They treat every task as if it were their first attempt, allocating the same computational resources at each step and repeating the same mistakes. This inability to accumulate knowledge from past experience significantly limits both their effectiveness and efficiency on complex, long-horizon tasks, where subroutines such as searching or creating notes are frequently reused across different objectives.

To address these limitations, we propose **Mobile-Agent-E**, a **hierarchical agentic framework** capable of **self-evolution** through past experiences. Mobile-Agent-E explicitly separates highlevel planning—such as decomposing tasks into subgoals—from low-level action decisions like determining specific actions and their parameters (e.g., tap(x,y)), enabling multi-level reasoning, perception, and error recovery. Figure 1 shows a demo of Mobile-Agent-E on a challenging online shopping task. Mobile-Agent-E also features a novel self-evolution module, which includes a persistent long-term memory containing two types of critical knowledge: *Tips* and *Shortcuts*. This design draws inspiration from human cognitive science, where Tips are akin to the lessons encoded in episodic memory [27], which involves recalling specific past experiences and using them to inform future decisions, while Shortcuts resemble procedural knowledge that facilitates the efficient and often subconscious execution of well-practiced tasks [23, 1].

To address the limitation of existing mobile benchmarks, which mainly include goal-oriented tasks, we introduce **Mobile-Eval-E**, a new challenging benchmark focusing on complex, open-ended, real-world tasks. Mobile-Eval-E features more than twice the number of expected operations per task compared to previous benchmarks [29, 38, 28] and incorporating a significantly higher proportion of multi-app tasks. Accompanying the benchmark, we introduce Satisfaction Score, a new metric to address the challenge posed by real-world tasks that often lack a binary success flag or a ground truth trajectory. This evaluation offers a reliable measure of agent performance aligned with real-world human needs.

To conclude, our contributions are threefold:

- **Mobile-Eval-E Benchmark**: A shift in focus from goal-oriented tasks to complex, open-ended, real-world mobile tasks.
- **Mobile-Agent-E Framework**: A hierarchical agent architecture featuring multi-level reasoning and error recovery, achieving a 22.1% absolute improvement over prior state-of-the-art approaches.
- **Self-Evolution Module**: The first work to explore self-evolution in mobile agents, yielding a 6.5% improvement in performance and a 12.9% gain in efficiency. Comprehensive analysis provides further insights for future research.

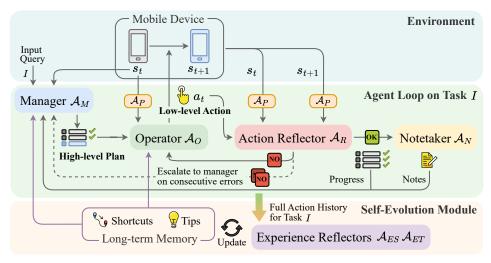


Figure 2: An overview of Mobile-Agent-E.

2 Mobile-Agent-E

Figure 2 provides an overview of Mobile-Agent-E. We detail the hierarchical agentic framework (§2.1) and the self-evolution module (§2.2) below.

2.1 Hierachical Agentic Framework

Multi-Level Reasoning. The key idea behind improving a model's performance on complex tasks is to disentangle high-level planning from low-level actions. Specifically, we divide all reasoning agents into two levels: the planning level, which contains a Manager, and the action level, which contains an Operator, an Action Reflector, and a Notetaker. Figure 3 provides a detailed breakdown of the inputs and outputs for each agent. All reasoning agents are instantiated from a frozen large multimodal model (LMM), such as GPT-40 [18]. We formally define each agent below, with notations given in Table 1.

The Manager (\mathcal{A}_M) focuses on devising high-level plans, i.e., identifying overall strategies and the next immediate subgoals, to fulfill the user's request. Note that the Shortcuts L_S (detailed in §2.2) are also provided to the Manager to guide efficient high-level planning.

$$W_P^t, W_S^t = \mathcal{A}_M(I, s_t, W_P^{t-1}, W_S^{t-1}, W_G^{t-1}, W_N^{t-1}, L_S)$$

Table 1: Summary of notations for intermediate inputs and outputs. Notations of agents are defined in §2.1.

Notation	Description
Environme	ent
I	Input task query
a^t	Action at step t
s^t	Phone state (screenshot)
Working M	1emory
W_V^t	Fine-grained visual perception
W_P^t	Overall plan (subgoals)
W_S^t	Current subgoal
$W_G^{\tilde{t}}$	Progress status
$W_N^{\overline{t}}$	Important notes
W_{EF}^t	Error Escalation Flag
$\mathbf{W}_{\mathbf{A}}$	Action history with outcome status
$\mathbf{W_E}$	Error history with feedback
Long-term	ı Memory
L_S	Shortcuts
L_T	Tips

The Operator (\mathcal{A}_O) decides which concrete action to perform based on the high-level plans from the Manager and the latest m-step history.* The Operator also considers the Tips as guidance from the long-term memory, which can be self-evolved from past experiences. Unlike the Manager, the action level agents take the fine-grained perception results W_V^t —in addition to the screenshot s_t —as input. We detail the perception module later in this section.

$$a_t = \mathcal{A}_O(I, s_t, W_V^t, W_P^t, W_S^t, W_G^t, W_N^t, \mathbf{W}_{\mathbf{A}}[-m:], \mathbf{W}_{\mathbf{E}}[-m:], L_S, L_T)$$

The action space of a_t is defined to contain not only *Atomic Operations* but also *Shortcuts*, which can evolve through tasks. The atomic operations include Open_App, Tap, Swipe, Type, Enter, Switch_App, Back, Home, and Wait. The full descriptions of the atomic operations can be found in Table 9. We detail the definitions and examples of *Shortcuts* and *Tips* in §2.2.

^{*}We empirically set m=5 in our experiments.

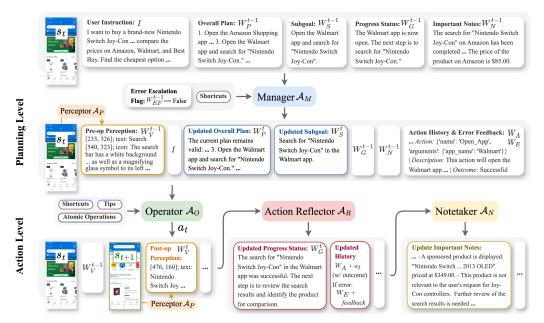


Figure 3: A detailed breakdown of one inference step t with Mobile-Agent-E, showing the inputs and outputs of each agent. Omitted information indicates no change.

The Action Reflector (A_R) checks the screenshots before (s_t) and after (s_{t+1}) of an action (a_t) to verify if the previous action achieves the expected outcome. We define three types of outcomes for an action: A. Successful or partially successful: the result of the last action meets the expectation; B. Failed: the last action results in a wrong page; and C. Failed: the last action produces no changes. After identifying the outcome, if the outcome is A, the Action Reflector updates the action history $\mathbf{W}_{\mathbf{A}}[t]$ as well as the progress status W_G^t . If the outcome is B or C, the Action Reflector additionally provides a description of the error and suggests potential reasons and solutions in $\mathbf{W}_{\mathbf{E}}[t]$.

$$\begin{aligned} \mathbf{W_A}[t], \mathbf{W_E}[t], W_G^t &= \mathcal{A}_R(I, s_t, W_V^t, s_{t+1}, W_V^{t+1}, \\ a_t, W_S^t, W_G^{t-1}) \end{aligned}$$

In complex mobile tasks, we often need to keep track of important notes during exploration, such as the price of a product or the phone number of a restaurant. The Notetaker (A_N) is dedicated to extracting and aggregating task-relevant information W_N^t after each step.

$$W_N^t = \mathcal{A}_N(I, s_{t+1}, W_V^{t+1}, W_P^t, W_S^t, W_G^t, W_N^{t-1})$$

Multi-Level Perception. The reasoning agents at different levels require different granularities of perception of the current phone state. At the *planning level*, the Manager only needs a holistic visual context of the screen to decide on high-level plans; therefore, we provide only the screenshot s_t to the Manager. At the *action level*, however, the agents need precise coordinates of interactive elements to predict and verify actions. To address this, we introduce the Perceptor (A_P) , a purely vision-based perception module that does not rely on the underlying XML file, following [28]. The Perceptor consists of three components: an OCR model, an icon-grounding model, and an icon-captioning model. Given a screenshot s_t , the Perceptor produces a fine-grained list of texts and icons along with their corresponding coordinates W_V^t .

Multi-Level Error Recovery. The ability to recover from errors is particularly important for executing complex, long-horizon tasks. In addition to the two-level reasoning agents, we introduce a two-level error-recovery mechanism that operates at both the action and planning levels. When an error first occurs (as reported by the Action Reflector), the Operator attempts to address it at the action level, for example, by tapping a different location. If the model becomes stuck in an error loop, i.e., it observes k consecutive failed actions (e.g., k=2), a special Error Escalation Flag, W_{EF}^{t-1} , is

[†]Some actions may need multiple repetitions to fulfill the expectation, for example, swipe up to find reviews. Thus, we include partially successful as meeting the expectation.

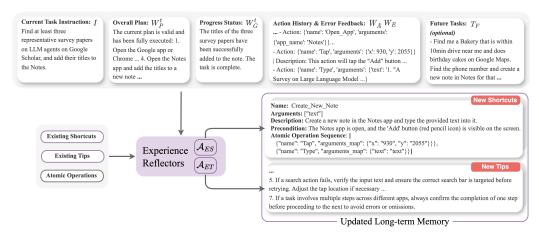


Figure 4: Illustration of the inputs and outputs to the Experience Reflectors for a single self-evolution step, including a concrete example of the newly generated Shortcuts and Tips.

raised to the Manager. In this case, the Manager receives additional information about the recent errors, $\mathbf{W_E}[-k:]$, and is asked to determine how to address the issue from a higher-level perspective, such as refining the overall plan or adjusting the current subgoal. A concrete example of how error escalation aids recovery is shown in Figure 9.

2.2 Self-Evolution Module

Inspired by how humans quickly adapt to new tasks, we maintain a long-term memory that persists across tasks and leverage two dedicated agents to reflect on past experiences. The long-term memory contains two important types of knowledge to evolve upon, **Tips** and **Shortcuts**, aiming to improve both the performance and efficiency of the model.

Tips (L_T) are defined as general guidance on effective interactions and lessons learned from previous trial-and-error experiences. Tips resemble episodic memory [27], which enables humans to recall past experiences and apply insights to future decisions.

Shortcuts (L_S) are defined as reusable, executable functions composed of sequences of atomic operations tailored for recurring subroutines. Shortcuts are akin to procedural knowledge, which allows humans to perform well-practiced tasks efficiently and often subconsciously [23, 1]. Due to the highly dynamic nature of the mobile environment, a Shortcut may only be applicable in certain states. For instance, the "Tap_Type_and_Enter" Shortcut (Figure 1) is usable only when the current screen has a text input box. To address this, we explicitly include a **precondition** in the definition of a Shortcut and require the Operator to verify that the current state satisfies the precondition before using the Shortcut. The arguments of a Shortcut have a unique one-to-one mapping to the arguments of its atomic operations.

When the self-evolution module is enabled, we leverage two Experience Reflectors, \mathcal{A}_{ES} and \mathcal{A}_{ET} , to update the Tips and Shortcuts based on the interaction history and optionally a list of future tasks T_F . The Experience Reflectors are also instantiated from frozen LMMs/LLMs. **Figure 4** provides a detailed breakdown of one self-evolution step. Figures 12 and 13 shows a full list of generated Shortcuts and Tips by Mobile-Agent-E.

3 Mobile-Eval-E Benchmark

3.1 Towards Complex, Open-Ended, Real-World Tasks

Existing mobile benchmarks—based on either simulated environments [20, 5] or dynamic actual devices [29, 38, 28]—primarily focus on goal-oriented tasks. These tasks often have a ground-truth trajectory and a unique end-state. In real-world scenarios, however, such tasks are unrealistic and leave a gap between benchmarking performance and practical applications.

Table 3: Comparison with state-of-the-art models on *complex open-ended tasks* in Mobile-Eval-E. GPT-40 is used as the LMM backbone for all methods.

Model	Туре	Binary Success Rate (%) ↑	Satisfaction Score (%) ↑			
	Traditiona	l Evaluation With	out Evolution	ı		
AppAgent [38]	Single-Agent	0.0	25.2	60.7	-	96.0
Mobile-Agent-v1 [29]	Single-Agent	4.0	45.5	69.8	-	68.0
Mobile-Agent-v2 [28]	Multi-Agent	8.0	53.0	73.2	96.7	52.0
Mobile-Agent-E (ours)	Multi-Agent	44.0	75.1	85.9	97.4	32.0
Evaluation With Cross-Task Evolution						
Mobile-Agent-E + Evo (ours)	Multi-Agent	40.0	86.9	90.4	97.8	12.0

Table 4: Results on different LMM backbones, including GPT-40, Gemini, and Claude. Metrics defined in §3.2.

Model	Gemini-1.5-pro	Claude-3.5-Sonnet	
	BS↑ SS↑ AA↑ RA↑ TE↓	BS↑ SS↑ AA↑ RA↑ TE↓	$ BS\uparrow SS\uparrow AA\uparrow RA\uparrow TE\downarrow$
	Traditional Evaluation	n Without Evolution	
Mobile-Agent-v2 [28]	4.0 50.8 63.4 83.9 64.0	12.0 70.9 76.4 96.9 32.0	8.0 53.0 73.2 96.7 52.0
Mobile-Agent-E (ours)	20.0 70.9 74.3 91.3 48.0	32.0 75.5 91.1 99.1 12.0	44.0 75.1 85.9 97.4 32.0
	Evaluation With Cre	ss-Task Evolution	
Mobile-Agent-E + Evo (ours	20.0 71.2 77.4 89.6 48.0	40.0 83.0 91.4 99.7 12.0	40.0 86.9 90.4 97.8 12.0

To address this limitation, we propose Mobile-Eval-E benchmark, which emphasizes complex, open-ended, real-world tasks. Mobile-Eval-E is designed for dynamic evaluation on actual devices and comprises 25 manually crafted tasks spanning five realistic scenarios: "Restaurant Recommendation," "Information Searching," "Online Shopping," "What's Trending," and "Travel Planning." Mobile-Eval-E tasks are long-horizon and reasoning-intensive, often admitting multiple satisfactory trajectories and success

Table 2: Comparison with existing dynamic evaluation benchmarks on real devices. Mobile-Eval-E emphasizes complex tasks that require significantly more steps and a wider variety of apps.

Benchmark	#Tasks	#Multi-App Tasks	#Apps	Avg # Steps	Total # Steps
Mobile-Eval	33	3	10	5.55	183
Mobile-Eval-v2	44	4	10	5.57	245
AppAgent	45	0	9	6.31	284
Mobile-Eval-E	25	19	15	14.56	364

states. As shown in Table 2, Mobile-Eval-E significantly surpasses previous dynamic benchmarks in complexity, featuring more than $2\times$ the number of expected steps per task. Additionally, Mobile-Eval-E encompasses a broader range of apps, with 76% of tasks requiring interactions with multiple apps. In \$4, we demonstrate that this benchmark poses a substantial challenge for existing state-of-the-art models. The full set of task queries can be found in Appendix Table \$.

3.2 Fine-Grained Evaluation Metrics

Previous benchmarks typically employ a *Binary Success Rate (BS)* or a completion rate against a "ground-truth" trajectory to evaluate task completeness. However, the complexity and open-endedness of Mobile-Eval-E tasks pose unique challenges in faithfully assessing model performance. For example, many tasks, such as "Plan a one-day itinerary for Palo Alto," may involve exploration and information aggregation, where multiple reasonable solutions might exist. Thus, we seek to measure *human satisfaction* rather than exact matches to ground-truth states.

For each task, we manually write a list of rubrics (an example is shown in Figure 5(a)), containing both milestone steps (e.g., "Opened Tripadvisor") and exploration criteria (e.g., "Viewed multiple attractions"). We then introduce the **Satisfaction Score (SS)** as the number of fulfilled rubrics divided by the total number of rubrics, as judged by a human evaluator. We also include *Action Accuracy (AA)* and *Reflection Accuracy (RA)* to evaluate action-level performance, and a *Termination Error (TE)* rate to reflect the agent's robustness and error-recovery capability. Details about the termination modes can be found in Appendix A. To keep the human evaluation workload reasonable for this fine-grained analysis, we maintain a relatively small number of tasks in Mobile-Eval-E.

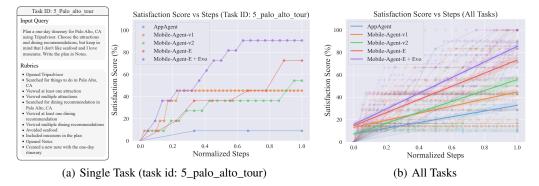


Figure 5: Satisfaction Score vs. Steps (SSS) curve for (a) a single task and (b) all tasks. In (a), we also present the human-written rubrics for the task. In (b), we additionally include a linear regression line for each model. To enable visualization of trajectories with different lengths on the same graph, we normalize the steps to the range [0, 1]. The y-axis of the rightmost point indicates the final satisfaction score. A steeper and higher line indicates better efficiency and effectiveness.

4 Experiments

We consider two evaluation settings—without and with evolution—to comprehensively assess both the hierarchical agentic framework and the self-evolution module. Our primary focus is on complex, reasoning-intensive tasks from Mobile-Eval-E. We also evaluate on Android World [20], which comprises 116 goal-oriented tasks across diverse apps. More details on baselines and model implementation can be found in Appendix A.

Traditional evaluation without evolution. In this setting, we evaluate each task individually without any cross-task information sharing. For Mobile-Eval-E, we perform dynamic, real-time evaluation [29, 38, 28] on a physical device. Specifically, we use the Android Debug Bridge (ADB) to control an Android phone[‡] and conduct human evaluation on the recorded screenshots and action histories. For Android World, we follow the official setup in an Android emulator[§], where evaluation is automated by verifying the final state. We adopt the same screen representation as M3A, including screenshots and the accessibility (A11y) tree.

Table 5: Comparison with state-of-the-art models on traditional *goal-oriented tasks* in Android World. We compare with methods released before Feb 2025.

Model	Success Rate (Pass@1 %)
Traditional Evaluation Without	Evolution
T3A + GPT-4-turbo [20]	30.6
M3A + GPT-4-turbo [20]	25.4
Ponder & Press + GPT-4o [31]	34.5
Aria-UI + GPT-4o [34]	44.8
UGround + GPT-4o [7]	44.0
Mobile-Agent-E + GPT-4o (ours)	45.0

Evaluation with cross-task evolution. To our knowledge, this is the first work to explore a cross-task evolution evaluation. In this setting, an agent is given a group of tasks and executes them sequentially, while maintaining a *persistent long-term memory* across tasks. Specifically, for Mobile-Eval-E, we form five groups corresponding to the five scenarios, each containing five tasks. We evaluate Mobile-Agent-E with the self-evolution module enabled (referred to as **Mobile-Agent-E + Evo**). At the end of the k-th task, the Experience Reflectors are prompted to update the long-term memory based on the interaction history of the current task as well as the queries for the remaining 5 - k tasks. This controlled setting allows us to investigate what is important for the agent to evolve from past experience (i.e., Tips and Shortcuts), and how this accumulated knowledge will impact subsequent tasks.

5 Results

Comparison with state-of-the-art. Tables 3 and 5 show that Mobile-Agent-E significantly outperforms prior SOTA (22.1%) on complex open-ended tasks, while also setting a new SOTA on traditional goal-oriented tasks in Android World. This comparison particularly highlights the effectiveness of the hierarchical agentic framework. Our approach also demonstrates superior robustness

[‡]We use Samsung Galaxy A15.

^{\$}https://github.com/google-research/android_world

Table 6: Analysis of computational overhead and Shortcut usage. In the inference speed table, the *reasoning only* section accounts for time spent solely on reasoning agents, while *perception* + *reasoning* includes the runtime of the Perceptor **on CPU**. Shortcut usage statistics are calculated as the ratio of Shortcuts used to the total number of actions performed by the Operator. The use of Shortcuts effectively accelerates inference, achieving comparable times to previous, simpler frameworks.

Inference Speed (Seconds per operation) \downarrow						
Model	Reas	oning O	nly	Percepti	ion + Red	asoning
Model	Gemini	Claude	GPT	Gemini	Claude	GPT
Mobile-Agent-v2	9.8	21.4	12.3	25.6	38.4	43.5
Mobile-Agent-E	16.5	25.5	17.4	30.8	41.0	30.1
Mobile-Agent-E + Evo	12.9	24.8	14.9	27.2	39.6	27.4

Shortcut Usage Percentage (%)					
Model	Gemini	Claude	GPT		
Mobile-Agent-E	11.9	12.8	12.4		
${\bf Mobile\text{-}Agent\text{-}E+Evo}$	14.8	13.2	14.4		

and error recovery capabilities, as indicated by a significantly lower Termination Error rate. Moreover, enabling self-evolution further enhances performance, leading to an improvement of 6.5% against no evolution. In §5.1, we provide further analysis of the evolution module.

Varying reasoning backbones. Table 4 demonstrates that Mobile-Agent-E can bring consistent improvements on all recent LMMs, including GPT-4o, Claude-3.5-Sonnet, and Gemini-1.5-pro. Moreover, we observe that self-evolution yields greater benefits when paired with stronger backbones.

Satisfaction Score v.s. Binary Success Rate. As shown in Tables 3 and 4, the Binary Success Rates (BS) are sparse and exhibit large jumps between different models. The Satisfaction Score (SS), on the other hand, provides a smoother, more faithful measure of the models' performance.

Task completion efficiency. Evaluating the efficiency of mobile agents on complex, open-ended tasks is not straightforward. Merely counting the number of steps is not optimal, as many tasks require exploration. A smaller number of steps reflects a quick exit but may result in insufficient exploration. Intuitively, if an agent fulfills more rubrics in a smaller number of steps, it is considered more efficient. Thus, we introduce the *Satisfaction Score vs Steps (SSS) curve* to compare and visualize the efficiency of different agents. To plot the SSS curve, we manually examine the recorded trajectories and track the satisfaction of rubrics after each step. As shown in Figure 5, we observe that Mobile-Agent-E not only achieves better final performance but also fulfills rubrics faster.

5.1 Further Analysis

Progressive impact of self-evolution over time.

The ideal behavior of self-evolution is to progressively bring more benefits to the agent as knowledge accumulates. To investigate this, we group the results of the tasks by their ordering index in each scenario and compare the performance with and without evolution. In Figure 6, the x-axis reflects the task index in the sequence it is performed. We observe a generally increasing trend indicating that the gain tends to be more significant in later tasks.

Shortcuts reduce computational overhead. The hierarchical multi-agent architecture in Mobile-Agent-E significantly improves performance on complex tasks but inevitably increases computational complexity. However, we found that the use of Shortcuts largely mitigates this overhead, enabling Mobile-Agent-E to achieve a speed comparable to that of

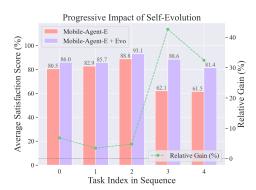


Figure 6: Later tasks in self-evolution show greater improvements, reflecting the growing impact of iterative evolution. The scores are averaged over five scenarios.

previous models. Note that when the self-evolution module is disabled, Mobile-Agent-E is provided with fixed initial Tips and a single example Shortcut. In Table 6, we observe a positive correlation between using more Shortcuts and faster inference speed. This is because a Shortcut enables the execution of multiple operations within a single decision-making iteration.

Unique impact from Tips. While the impact from Shortcuts is directly visible in the action history, it is less obvious whether the evolved Tips bring distinctive benefits.



Figure 7: Case study example where relevant Shortcuts and Tips are automatically retrieved from the previously evolved long-term memory and subsequently leveraged to complete an unseen, challenging task. The action trajectory also includes an example where the agent recovers from an error.

To ablate on this, we filter out task instances where the same set of unique Shortcuts is used or where only atomic actions are employed, and compare the Satisfaction Score with or without the evolved Tips. Table 7 shows that Tips alone serve as an important aspect of self-evolution.

Table 7: Unique impact from the evolved Tips.

	Gemini	Claude	GPT-40
Mobile-Agent-E	69.0	75.6	79.7
Mobile-Agent-E + evolved Tips	72.6	85.2	87.5

6 Managing Self-evolution in the Long Run

In real-world mobile usage, after running the agent on a large number of tasks in various scenarios, the accumulated Tips and Shortcuts may grow to an amount where it is no longer feasible to include all of them in the decision-making context. Thus, in this case study, we aim to explore closing the self-evolution loop by introducing two additional **Experience Retriever** agents for Tips \mathcal{A}_{ERT} and Shortcuts \mathcal{A}_{ERS} . We consider a new task in an unknown scenario, as shown in Figure 7. First, we provide all the updated Tips and Shortcuts—after running Mobile-Agent-E on all 5 scenarios (a total of 25 tasks) in Mobile-Eval-E—to the Experience Retrievers. With GPT-40 as the backbone, the updated long-term memory contains a total of 7 unique Shortcuts and 59 Tips, among which 6 Shortcuts and 55 Tips are newly proposed by Mobile-Agent-E during experience reflection. Then, the Experience Retrievers are prompted to select only the relevant Tips and Shortcuts for the current task. The qualitative example in Figure 7 shows that Mobile-Agent-E effectively retrieves and leverages highly relevant Shortcuts and Tips to successfully complete a challenging unseen task. The full list of Tips and Shortcuts after evolution can be found in Appendices H and G.

7 Related Work

7.1 GUI Agents

The advancement of large multimodal models (LMM) has driven research on LMM-based GUI agents [30], focusing on AI assistants for GUI environments like Web [6, 39, 8, 35, 21], PC [9, 37, 15, 33, 24], and mobile devices [29, 38, 11, 28, 14]. For mobile, research has enhanced single-agent perception and reasoning via tool usage [29] and exploration [38, 11]. Recent multi-agent systems [20, 28] show promise but still face challenges like short-sighted planning and poor error recovery. It is worth noting that the "planning" module in Mobile-Agent-v2 [28] merely serves as a progress tracker and is *fundamentally different* from the proposed Manager in Mobile-Agent-E. In Mobile-Agent-v2, the "decision-making" module remains responsible for both high-level planning (e.g., "what to do next") and low-level action execution (e.g., "where to tap"), resulting in a flat

and overloaded design. In contrast, Mobile-Agent-E introduces a hierarchical agentic structure that explicitly separates high-level planning from low-level action decisions.

7.2 Self-Evolution in Foundation Models

Self-improvement in large language and multimodal models has been widely explored [25], through techniques like iterative refinement [16], self-reflection [22], self-training [10], and multi-persona collaboration [32]. Recent work also emphasizes tool learning and creation [4, 19, 36]. In GUI agents, self-evolution is less explored. While Cradle [24] shows potential in skill curation for PC environments, how to do evolution in mobile settings remains unaddressed. In this work, we identify two important types of knowledge for evolution, namely Tips and Shortcuts.

8 Conclusion and Limitations

In this work, we make the first attempt to build mobile agents for complex, real-world tasks, demonstrating the effectiveness of hierarachical agentic framework as well as self-evolution. Future work will focus on developing improved strategies for generating, invoking, and revising long-term memory, as well as automating the evaluation of complex mobile tasks. As detailed below, remaining limitations include erroneous perception, imperfect Shortcut generation, and human evaluation.

Misuse of Shortcuts due to incorrect perception of phone state. Although we explicitly require the Operator to verify the current phone state to ensure it fulfills the *precondition* of a Shortcut before calling it, there are still cases where the model incorrectly perceives the state, resulting in the misuse of Shortcuts in an invalid state. Figure 10 illustrates an example of such error. A detailed description of the example is provided in the caption. This type of error could potentially be mitigated by employing a dedicated agent for verifying preconditions or by enhancing the perception module to better understand phone states.

Errors and imperfections in self-evolved shortcuts. Although effective in most cases, we still observe errors and imperfections in the agent-generated Shortcuts during self-evolution. These issues can lead to propagated errors when an erroneous Shortcut is used in subsequent tasks. Figure 11 illustrates an example of such erroneous and imperfect Shortcuts. A detailed description of the example is provided in the caption. This highlights the need for future work on approaches to generate higher-quality Shortcuts and equipping the agent with the ability to reflect on and revise generated Shortcuts in subsequent tasks.

Human evaluation is required for Mobile-Eval-E. Due to the complexity and open-ended nature of these real-world tasks, no current multimodal model can replace human evaluation for providing a reliable assessment. In the future, as multimodal reasoning models advance, we believe it is possible to develop an automatic rubric verifier that determines whether each criterion is met at a given step based on the global action history. We leave this promising direction to future work.

Broader Impacts

This paper aims to advance the field of LMM-based agents by developing a hierarchical multiagent framework and benchmark to improve the usability and efficiency of smartphones in complex, multi-step tasks.

While the primary aim is to enhance mobile task efficiency and user accessibility, the development of mobile agents capable of autonomous decision-making introduces potential risks. For example, unauthorized or unintended actions by the agent, such as the misuse of sensitive information including credit card details or private data, could result in serious consequences for users. These risks emphasize the critical need for robust safeguards, error recovery mechanisms, and fail-safe systems to ensure that the agent's actions consistently align with user intentions.

We are actively pursuing future work that focuses on designing and integrating robust privacy and safety mechanisms. These include explicit user consent workflows for sensitive operations, encryption protocols to protect user data during processing and storage, and automated systems to flag potentially harmful or unauthorized actions. These advancements will be crucial for maximizing the societal benefits of these systems and building user trust in autonomous mobile agents.

References

- [1] John R Anderson. Acquisition of cognitive skill. Psychological review, 89(4):369, 1982.
- [2] Anthropic. Claude 3.5 Sonnet, 2024.
- [3] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-VL: A Frontier Large Vision-Language Model with Versatile Abilities. *arXiv* preprint arXiv:2308.12966, 2023.
- [4] Tianle Cai, Xuezhi Wang, Tengyu Ma, Xinyun Chen, and Denny Zhou. Large language models as tool makers. *arXiv preprint arXiv:2305.17126*, 2023.
- [5] Jingxuan Chen, Derek Yuen, Bin Xie, Yuhao Yang, Gongwei Chen, Zhihao Wu, Li Yixing, Xurui Zhou, Weiwen Liu, Shuai Wang, et al. Spa-bench: A comprehensive benchmark for smartphone agent evaluation. In NeurIPS 2024 Workshop on Open-World Agents, 2024.
- [6] Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and Yu Su. Mind2web: Towards a generalist agent for the web. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- [7] Boyu Gou, Ruohan Wang, Boyuan Zheng, Yanan Xie, Cheng Chang, Yiheng Shu, Huan Sun, and Yu Su. Navigating the digital world as humans do: Universal visual grounding for gui agents. *arXiv preprint arXiv:2410.05243*, 2025.
- [8] Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Yong Dai, Hongming Zhang, Zhenzhong Lan, and Dong Yu. Webvoyager: Building an end-to-end web agent with large multimodal models. *arXiv preprint arXiv:2401.13919*, 2024.
- [9] Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxiao Dong, Ming Ding, and Jie Tang. Cogagent: A visual language model for gui agents, 2023.
- [10] Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. Large language models can self-improve. *arXiv preprint arXiv:2210.11610*, 2022.
- [11] Yanda Li, Chi Zhang, Wanqi Yang, Bin Fu, Pei Cheng, Xin Chen, Ling Chen, and Yunchao Wei. Appagent v2: Advanced agent for flexible mobile interactions. *arXiv preprint arXiv:2408.11824*, 2024.
- [12] Minghui Liao, Zhaoyi Wan, Cong Yao, Kai Chen, and Xiang Bai. Real-time scene text detection with differentiable binarization. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 11474–11481, 2020.
- [13] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, and Lei Zhang. Grounding DINO: marrying DINO with grounded pre-training for open-set object detection. CoRR, abs/2303.05499, 2023.
- [14] Xiao Liu, Bo Qin, Dongzhu Liang, Guang Dong, Hanyu Lai, Hanchen Zhang, Hanlin Zhao, Iat Long Iong, Jiadai Sun, Jiaqi Wang, et al. Autoglm: Autonomous foundation agents for guis. *arXiv preprint arXiv:2411.00820*, 2024.
- [15] Xiao Liu, Tianjie Zhang, Yu Gu, Iat Long Iong, Yifan Xu, Xixuan Song, Shudan Zhang, Hanyu Lai, Xinyi Liu, Hanlin Zhao, et al. Visualagentbench: Towards large multimodal models as visual foundation agents. *arXiv preprint arXiv:2408.06327*, 2024.
- [16] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. Advances in Neural Information Processing Systems, 36, 2024.
- [17] Dang Nguyen, Jian Chen, Yu Wang, Gang Wu, Namyong Park, Zhengmian Hu, Hanjia Lyu, Junda Wu, Ryan Aponte, Yu Xia, et al. Gui agents: A survey. *arXiv preprint arXiv:2412.13501*, 2024.
- [18] OpenAI. GPT-4o System Card, 2024.
- [19] Cheng Qian, Chi Han, Yi R Fung, Yujia Qin, Zhiyuan Liu, and Heng Ji. Creator: Tool creation for disentangling abstract and concrete reasoning of large language models. arXiv preprint arXiv:2305.14318, 2023.
- [20] Christopher Rawles, Sarah Clinckemaillie, Yifan Chang, Jonathan Waltz, Gabrielle Lau, Marybeth Fair, Alice Li, William Bishop, Wei Li, Folawiyo Campbell-Ajala, et al. Androidworld: A dynamic benchmarking environment for autonomous agents. arXiv preprint arXiv:2405.14573, 2024.

- [21] Revanth Gangi Reddy, Sagnik Mukherjee, Jeonghwan Kim, Zhenhailong Wang, Dilek Hakkani-Tur, and Heng Ji. Infogent: An agent-based framework for web information aggregation. *arXiv* preprint *arXiv*:2410.19054, 2024.
- [22] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. Advances in Neural Information Processing Systems, 36, 2024.
- [23] Larry R Squire and Stuart M Zola. Structure and function of declarative and nondeclarative memory systems. Proceedings of the National Academy of Sciences, 93(24):13515–13522, 1996.
- [24] Weihao Tan, Ziluo Ding, Wentao Zhang, Boyu Li, Bohan Zhou, Junpeng Yue, Haochong Xia, Jiechuan Jiang, Longtao Zheng, Xinrun Xu, et al. Towards general computer control: A multimodal agent for red dead redemption ii as a case study. In ICLR 2024 Workshop on Large Language Model (LLM) Agents, 2024.
- [25] Zhengwei Tao, Ting-En Lin, Xiancai Chen, Hangyu Li, Yuchuan Wu, Yongbin Li, Zhi Jin, Fei Huang, Dacheng Tao, and Jingren Zhou. A survey on self-evolution of large language models. arXiv preprint arXiv:2404.14387, 2024.
- [26] Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- [27] Endel Tulving. Episodic memory: From mind to brain. Annual review of psychology, 53(1):1–25, 2002.
- [28] Junyang Wang, Haiyang Xu, Haitao Jia, Xi Zhang, Ming Yan, Weizhou Shen, Ji Zhang, Fei Huang, and Jitao Sang. Mobile-agent-v2: Mobile device operation assistant with effective navigation via multi-agent collaboration. *arXiv* preprint arXiv:2406.01014, 2024.
- [29] Junyang Wang, Haiyang Xu, Jiabo Ye, Ming Yan, Weizhou Shen, Ji Zhang, Fei Huang, and Jitao Sang. Mobile-agent: Autonomous multi-modal mobile device agent with visual perception. arXiv preprint arXiv:2401.16158, 2024.
- [30] Shuai Wang, Weiwen Liu, Jingxuan Chen, Weinan Gan, Xingshan Zeng, Shuai Yu, Xinlong Hao, Kun Shao, Yasheng Wang, and Ruiming Tang. Gui agents with foundation models: A comprehensive survey. arXiv preprint arXiv:2411.04890, 2024.
- [31] Yiqin Wang, Haoji Zhang, Jingqi Tian, and Yansong Tang. Ponder & press: Advancing visual gui agent towards general computer control. *arXiv preprint arXiv:2412.01268*, 2024.
- [32] Zhenhailong Wang, Shaoguang Mao, Wenshan Wu, Tao Ge, Furu Wei, and Heng Ji. Unleashing the emergent cognitive synergy in large language models: A task-solving agent through multi-persona selfcollaboration. arXiv preprint arXiv:2307.05300, 2023.
- [33] Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, et al. Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments. *arXiv preprint arXiv:2404.07972*, 2024.
- [34] Yuhao Yang, Yue Wang, Dongxu Li, Ziyang Luo, Bei Chen, Chao Huang, and Junnan Li. Aria-ui: Visual grounding for gui instructions. *arXiv preprint arXiv:2412.16256*, 2024.
- [35] Ori Yoran, Samuel Joseph Amouyal, Chaitanya Malaviya, Ben Bogin, Ofir Press, and Jonathan Berant. Assistantbench: Can web agents solve realistic and time-consuming tasks?, 2024.
- [36] Lifan Yuan, Yangyi Chen, Xingyao Wang, Yi R Fung, Hao Peng, and Heng Ji. Craft: Customizing Ilms by creating and retrieving from specialized toolsets. *arXiv preprint arXiv:2309.17428*, 2023.
- [37] Chaoyun Zhang, Liqun Li, Shilin He, Xu Zhang, Bo Qiao, Si Qin, Minghua Ma, Yu Kang, Qingwei Lin, Saravan Rajmohan, Dongmei Zhang, and Qi Zhang. UFO: A UI-Focused Agent for Windows OS Interaction. arXiv preprint arXiv:2402.07939, 2024.
- [38] Chi Zhang, Zhao Yang, Jiaxuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. Appagent: Multimodal agents as smartphone users, 2023.
- [39] Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. Gpt-4v(ision) is a generalist web agent, if grounded. In *Forty-first International Conference on Machine Learning*, 2024.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The claims are clearly presented in the abstract and intro, and well supported by empirical results and analysis in the Sections 4 and 5.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Limitations are discussed in Section 7 and visualized in Appendix E.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: NA Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Implementation details are presented in Appendix B.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The data is open-sourced, the instruction can be found in supplementary Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: All experimental settings are detailed in Section 3 and Appendix B.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [No]

Justification: We did not provide error bars due to high cost of repetitive runs using APIs. Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We provide the required APIs in Appendix B and discussed time of execution in Section 5.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: The paper conform, in every respect, with the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: We will add the broader impacts in the additional page once accepted.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]
Justification: NA

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All code, datasets and models used are well cited.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

• If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: All new assets are well documented in the supplementary

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [Yes]

Justification: We include a detailed dicussion on manual efforts in benchmark creation and evaluation. We also include the evaluation script in the supplementary.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]
Justification: NA

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]
Justification: NA
Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

A Experimental Details

Baselines. For Mobile-Eval-E, we compare against diverse open-sourced mobile agent frameworks compatible with actual devices, including AppAgent [38], Mobile-Agent-v1 [29], and Mobile-Agent-v2 [28]. To maximize an apple-to-apple comparison with Mobile-Agent-v2, which is the previous state-of-the-art, we apply an identical atomic operation space, perception model, and initial Tips to Mobile-Agent-v2 as Mobile-Agent-E. AppAgent originally requires an additional exploration phase, which does not fit our setting; thus, we add the initial Tips as additional knowledge. AppAgent-v2 [11] is not included since it is not open-sourced at the time of writing.

We also compare with a wider range of models with reported scores on Android World, such as M3A [20], Aria-UI [34] and UGround [7]. We further explore using different large multimodal models (LMM) as backbones for the reasoning agents, including GPT-40 [18], Claude-3.5-Sonnet [2], and Gemini-1.5-pro [26]. Unless otherwise specified, the default backbone for all models is GPT-40.

Backbone versions. The detailed versions of the large multimodal models are listed as follows: (1) GPT-40 version: gpt-40-2024-11-20; (2) Claude-3.5 version: claude-3-5-sonnet-20241022; (3) Gemini-1.5 version: gemini-1.5-pro-latest (Dec 2024)

Perceptor implementation details. We follow Mobile-Agent-v2 [28] to implement the Perceptor in Mobile-Agent-E with slight modifications. We use DBNet[¶][12] and ConvNextViT-document[¶] from ModelScope for OCR detection and recognition respectively. We use GroundingDINO [13] for icon grounding and Qwen-VL-Plus [3] for generating captions for each cropped icon.

Agent termination modes. There are five ways an agent can exit from performing a task: (1) self-reported success: the agent decides to stop on its own; (2) reaching the maximum number of iterations: we set the maximum iteration count to 40 to prevent infinite loops; (3) reaching the maximum number of consecutive errors: if the agent has an action reflector and it identifies 3 consecutive errors, the agent is exited; (4) reaching the maximum number of repeated actions: if the agent performs the exact same action (excluding Swipe and Back) more than 3 consecutive times; (5) any other errors, such as errors when parsing the raw response into a valid action. If a task exits in one of the ways described in 2–5, it is marked as having a termination error (TE). The TE rate is computed as the ratio of tasks with termination errors to all tasks.

Note that we recognize self-reported success may be inaccurate; however, in the TE metric we treat it as a normal termination. Early termination is already penalized via a low satisfaction score, while the TE metric is designed to capture the model's robustness to truly abnormal terminations—such as infinite loops, formatting errors, and so on.

B Full Trajectory Comparison Example with Previous SOTA

Figure 8 presents the full trajectory of the task shown in Figure 1, comparing the previous state-of-the-art, Mobile-Agent-v2 [28], and our proposed Mobile-Agent-E. Mobile-Agent-v2 suffers from early termination after interacting with two Apps, whereas Mobile-Agent-E fulfills all rubrics and stops at the App offering the best deal.

C Impact of Multi-level Error Recovery

Figure 9 illustrates how the error escalation mechanism in Mobile-Agent-E enhances error recovery ability. A detailed description of the example is provided in the caption.

D Illustration of Remaining Limitations

Figure 10 illustrates an example of a misuse of Shortcuts in an invalid state. Figure 11 illustrates an example of erroneous and imperfect Shortcuts.

 $[\]P$ https://modelscope.cn/models/iic/cv_resnet18_ocr-detection-db-line-level_damo

https://modelscope.cn/models/iic/cv_convnextTiny_ocr-recognition-document_damo



Figure 8: Full trajectory comparison between the previous state-of-the-art, Mobile-Agent-v2 [28], and Mobile-Agent-E.

E All Tasks in Mobile-Eval-E Benchmark

Table 8 presents the input queries, involved App types, and scenarios for all Mobile-Eval-E tasks. The complete list of rubrics and human reference operation sequences is provided in the supplementary material.

F Atomic Operation Space

Table 9 presents all atomic operations considered in Mobile-Agent-E.

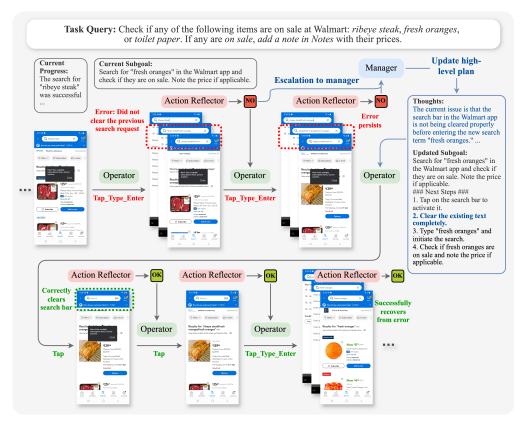


Figure 9: Error recovery with escalation. The task requires the agent to search for three different items on Walmart and note their sales information. At the step shown in the figure, the agent has already searched for ribeye steak and intends to search for fresh oranges next. However, the Operator erroneously calls the Shortcut that inputs text into the search bar and performs a search without clearing the previously entered text. Although the Action Reflector raises an error, the subgoal remains unchanged, and the Operator fails to rectify the error on the second attempt. After observing two consecutive errors, the error is escalated to the Manager, which correctly identifies the problem and revises the subgoal with detailed, decomposed steps to address the error. This helps the Operator correctly recover from the previous error by first tapping the "×" icon to clear the previous search query.

G Full list of Self-Evolved Shortcuts

Figure 12 shows a full list of generated Shortcuts by Mobile-Agent-E after self-evolution on all 25 tasks from Mobile-Eval-E benchmark.

H Full list of Self-Evolved Tips

Figure 13 shows a full list of generated Tips by Mobile-Agent-E after self-evolution on all 25 tasks from Mobile-Eval-E benchmark.

I Manual Efforts in Benchmark Construction and Evaluation

Task Curation: We first curate the scenarios and tasks via a collaboration with OpenAI GPT-4o. We (the authors) first provide a set of seed scenarios and task examples, and GPT-4o instantiates. We then manually examine and curate the tasks to ensure (1) feasibility; (2) diversity; (3) complex enough; (4) no ambiguity; (5) having some shared subgoals.

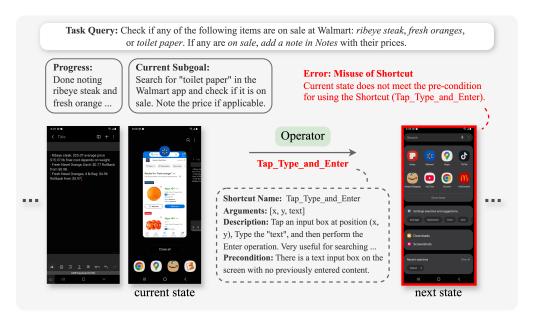


Figure 10: Example of misuse of Shortcuts in an invalid state. At the current step, as shown in the figure, the agent intended to switch back to Walmart to search for the final item requested by the user. While it correctly performs the "Switch_App" operation, it then calls a Shortcut for searching without realizing that it is not yet in the App where the search bar is available.

Rubrics Curation: We first ask one of the authors to manually perform each task on the device, logging the operation steps (which is used as a reference in Table 2); Then we ask the author to write a list of rubrics for each task following the criteria: (1) the rubrics should be objective; (2) the rubrics should cover key milestones; (3) the rubrics should give credit to reasonable exploration, e.g., "checked at least two reviews". Finally, we do cross examination among all the authors to ensure the rubrics are of good quality.

Evaluation: The authors are responsible for performing manual evaluation on the recorded trajectories for Mobile-Eval-E tasks. The evaluation tool is included in the supplementary code. It is important to note that, during evaluation, we do not enforce an order for fulfilling the rubrics. For example, the model can fulfill "checked the price of ribeye" and "checked the price of fresh oranges" with any order in the task "3_walmart_sale_items".

Figure 11: Example of imperfect (above) and erroneous (below) generated Shortcuts. The "Search_Location_in_Maps" Shortcut includes an unnecessary Tap action in the operation sequence, while the "Switch_App_And_Search" Shortcut omits a Tap action needed to first enter the desired App before performing the search.

 Table 8: All task queries in Mobile-Eval-E.

Scenario	Task ID	APPs	Input Query
	1_late_night_korean_food	Maps	$Find the best-rated \ late-night Korean \ restaurant \ in \ Champaign, \ IL \ that \ opens \ beyond \ 9pm \ on \ Google \ Maps.$
Restaurant Recommendation	1_nearest_bakery	Maps	Get directions to the nearest Bakery that has a rating higher than 4.0 on Google Maps. Stop at the screen showing the route.
	1_thai_duck	Maps, Notes	Find the best-rated Thai restaurant in Urbana, IL that serves duck cuisine on Google Maps. Review customer comments and compile a summary of positive and negative feedback in Notes.
	1_bakery_birthday_cake	Maps, Notes	Find me a Bakery that is within 10min drive near me and does birthday cakes on Google Maps. Find the phone number and create a new note in Notes for that.
	1_chinese_ohare	Maps, X, Notes	Find me a popular Chinese restaurant near Chicago O'Hare airport on Google Maps. Check X for recent posts about their signature dishes and write a summary in Notes. Then get directions to that restaurant on Google Maps. Stop at the screen showing the route.
	2_segment_anything_cited	Chrome	Find the most-cited paper that cites the paper 'Segment Anything' on Google Scholar. Stop at the screen showing the paper abstract.
	2_llm_agents_survey	Chrome, Notes	Find at least three representative survey papers on LLM agents on Google Scholar, and add their titles to the Notes.
Information Researching	2_recipes_chinese	Chrome, YouTube	I have some onions, beef, and potatoes in my refrigerator. Can you find me a Chinese-style recipe that uses all three ingredients and can be prepared in under an hour? And find me a video tutorial or YouTube for that. Stop at the screen displaying the video.
	2_mcdonalds_deals	McDonald's, Maps	Can you check the McDonald's APP to see if there are any Rewards or Deals including Spicy McCrispy. If so, help me add that to Mobile Order (Do not pay yet, I will do it myself). And ther check the pickup location and get directions on Google Maps. Stop at the screen showing the route
	2_headphones_reviews	Amazon, Notes	Find three detailed user reviews of the Bose QC45 headphones from Amazon. Summarize the general sentiment in the Notes.
Online Shopping	3_oled_tv	Best Buy	Find the best deal on a 55 -inch $4K$ OLED TV at Best Buy. Stop at the screen displaying the best deal you find.
	3_laptop_nvidia_gpu	Amazon Shop- ping	Find me a laptop on Amazon that is under \$1000 with an Nvidia GPU and more than 8GB RAM.
	3_ninja_air_fryer	Amazon Shop- ping, Walmart	Compare the price of a Ninja air fryer 8 qt at Walmart and Amazon. Stop at the screen displaying the best deal you find.
	3_walmart_sale_items	Walmart, Notes	Check if any of the following items are on sale at Walmart: ribeye steak, fresh oranges, or toile paper. If any are on sale, add a note in Notes with their prices.
	3_nintendo_switch_joy_con	Amazon Shop- ping, Best Buy, Walmart	I want to buy a brand-new Nintendo Switch Joy-Con. Any color is fine. Please compare the prices on Amazon, Walmart, and Best Buy. Find the cheapest option and stop at the screen where I can add it to the cart.
	4_x_black_myth_wukong	X, Notes	Find the top posts about the game 'Black Myth Wukong' on X and summarize the key highlights in Notes.
	4_x_trending_news	X, Notes	Check the top 3 trending news on X. Read a few posts to figure out what's happening. And create a new Note to summarize your findings.
What's Trending	4_watercolor_painting_tutoria	alLemon8, Notes	I want to learn how to paint watercolor. Find me some content creators to follow on Lemon8 that has highly liked posts about watercolor painting tutorials. List their account names in Notes.
	4_movie_trending	Fandango, Notes	Check the top 5 trending movies on Fandango that are currently in theaters. Compare their ratings and create a note in Notes for the highest-rated one, including its name and showtimes.
	4_horror_movie_reviews	Fandango, Lemon8, Notes	Find me the latest horror movie currently in theaters on Fandango. Check some reviews on Lemont about the movie and create a note in Notes with the general sentiment.
	5_cheap_flights_newyork	Booking	Find the cheapest round-trip flight from Chicago to New York City in the next month on Booking Stop at the screen showing the best deal.
	5_things_to_do_la	Tripadvisor, Notes	Suggest some interesting things to do in LA. Find the top 3 attractions on Tripadvisor. Save the lis in Notes.
Travel Planning	5_palo_alto_tour	Tripadvisor, Notes	Plan a one-day itinerary for Palo Alto, CA using Tripadvisor. Choose the attractions and dining recommendations, but keep in mind that I don't like seafood and I love museums. Write the plan in Notes.
	5_local_food_chicago	Tripadvisor, Notes	Find a highly recommended local restaurant in Chicago on Tripadvisor. Check the reviews about must-try dishes and summarize in Notes.
	5_hotel_champaign	Booking, Maps	Help me find a hotel in Champaign, IL on Booking that is under \$200 for a queen bed. Make sure that the rating is higher than 7.0. Double-check on Google Maps to see if it is close to Green Street

Table 9: Atomic operations space.

Operation	Description
$\overline{Open_App(app_name)}$	If the current screen is Home or App screen, you can use this action to open the app named "app_name" on the visible on the current screen.
Tap(x,y)	Tap the position (x, y) in current screen.
$Swipe(x_1, y_1, x_2, y_2)$	Swipe from position (x_1,y_1) to position (x_2,y_2) . To swipe up or down to review more content, you can adjust the y-coordinate offset based on the desired scroll distance. For example, setting $x_1=x_2=0.5*width,\ y_1=0.5*height,\ and\ y_2=0.1*height\ will swipe upwards to review additional content below. To swipe left or right in the App switcher screen to choose between open apps, set the x-coordinate offset to at least 0.5*width.$
Type(text)	Type the "text" in an input box.
Enter()	Press the Enter key after typing (useful for searching).
$Switch_App()$	Show the App switcher for switching between opened apps.
Back()	Return to the previous state.
Home()	Return to home page.
Wait()	Wait for 10 seconds to give more time for a page loading.

Figure 12: Full list of Shortcuts generated by Mobile-Agent-E (with GPT-40) after self-evolution.

** Initial Tips (User Provided) **

0. Do not add any payment information. If you are asked to sign in, ignore it or sign in as a guest if possible. Close any pop-up windows when opening an app. 1. By default, no apps are opened in the background. 2. Screenshots may show partial text in text boxes from your previous input; this does not count as an error. 3. When creating new Notes, you do not need to enter a title unless the user specifically requests it.

** Agent Generated Tips (Scenario 1) **

**A. When searching for restaurants or businesses, ensure the query includes specific details like location, type of cuisine, and operational hours to narrow down results effectively. 5.

Always verify the operational hours of businesses to ensure they meet the user's requirements, especially for late-night or time-sensitive searches. 6. When filtering search results (e.g., by rating or distance), ensure the filter criteria are applied correctly to avoid irrelevant results. 7. Double-check the selected location or business to ensure it matches the user's requirements (e.g., rating, proximity, or specific services offered) before proceeding to the route screen. 8. If the task index certaing a route, confirm that the route is displayed correctly and matches the intended destination before marking the subgoal as complete. 9. When navigating through menus or categories, use a systematic approach to ensure all relevant sections are explored thoroughly. 10. If an action does not return to the expected screen, use alternative navigation methods (e.g., tapping "X" or returning to the home screen) to correct the workflow. 11.

When summarizing customer feedback, include both positive and negative aspects to provide a balanced overview. 12. When retrieving contact information, ensure the details (e.g., phone number or address) are accurate and match the selected business before saving them in Notes. 13. If a task involves multiple apps (e.g., Google Maps and Notes), ensure smooth transitions between apps and verify that the required information is correctly transferred. 14. If an app fails to open or respond in the app switcher, return to the home screen and reopen the app directly to avoid delays.

** Agent Generated Tips (Scenario 2) **

A When identifying the most-cited paper or similar tasks, ensure to sort the results by citation count if the option is available. This minimizes manual scanning and reduces errors. 5. If a search action fails, verify the input text and ensure the correct search bar is targeted before retrying. Adjust the tap location if necessary. 6. When recording information from search results, ensure the details are accurate and clearly formatted to avoid confusion. 7. If a task involves multiple steps across different apps, always confirm the completion of one step before proceeding to the next to avoid errors or omissions. 8. If a search query fails to execute, double-chebt the tap location and ensure the search bar is properly activated before typing. Retry the action with slight adjustments if necessary. 9. When selecting a video or item from a list, ensure the title matches the intended choice to avoid selecting the wrong option. 10. If a button or option does not respond to a tap, ensure it is fully visible on the screen. Use a swipe or scroll action to adjust the view if necessary before retrying. 11. When switching between apps, ensure the correct app is selected from the app switcher to avoid unnecessary navigation errors. 12. Always stop at the final screen requested by the user, ensuring the task is fully completed before ending the interaction.

4. When identifying the best deal, prioritize both price and features, and ensure any discounts or promotions are clearly noted. 5. Always confirm that the displayed product matches the search criteria (e.g., size, specifications) to avoid selecting an incorrect item. 6. If the task requires stopping at a specific screen, ensure the screen is fully loaded and all relevant details are visible before stopping. 7. If a filter does not apply correctly, try adjusting it again by swiping or tapping alternative areas of the screen to reveal hidden options. 8. When using sliders for filters (e.g., price range), swiping is often more effective than tapping to adjust the values. 9. If a filter unexpectedly resets or removes itself, reapply it and verify the results before proceeding. 10. Always double-check the final results to ensure all filters (e.g., price, specifications) have been applied correctly. 11. When comparing prices across platforms, ensure that the product model and specifications (e.g., size, features) are identical to avoid inaccurate comparisons. 12. If swiping to reveal content, ensure the swipe is smooth and covers enough distance to load all relevant details on the screen. 13. If an app fails to open or navigate correctly, return to the home screen and retry the action. This often resolves navigation issues. 14. If a tap action does not work as expected, consider tapping alternative areas of the screen, such as associated buttons or options, to achieve the desired outcome. 15. When switching between apps, ensure the correct app is reopened and verify the screen before proceeding to avoid unnecessary repetition.

** Agent Generated Tips (Scenario 4) **

4. When navigating apps, ensure that the correct icon is tapped by carefully identifying its position and function to avoid misalignment or unintended actions. 5. If a search filter is applied unintentionally, clear it by tapping the "X" icon in the search bar before proceeding with a new search. 6. Always verify the context of the search results to ensure they align with the intended query before summarizing or proceeding to the next step. 7. When recording information in cless, ensure the formatting is clear and consistent for easy readability. 8. Double-check the accuracy of the recorded information (e.g., account names, titles) before saving the note to avoid errors. 9. If redirected to an unintended page (e.g., "My Orders"), an avigate back to the main interface or intended section before proceeding. 10. When comparing multiple items (e.g., mover ratings), keep track of all relevant data to ensure accurate comparisons and avoid revisiting the same pages unnecessarily. 11. If an app opens an unintended interface (e.g., camera instead of Notes), return to the home screen and retry opening the correct app to avoid confusion. 12. When entering search terms, ensure the previous query is cleared completely to prevent appending incorrect text to the new query. 13. If a misaligned tap opens an unintended menu (e.g., Filters), close it immediately and retry the intended action. 14. Use broader search terms if specific queries fail to yield results, and refine the search gradually based on the context. 15. If an app fails to execute a search or action, consider switching to a browser or alternative app to complete the task.

** Agent Generated Tips (Scenario 5) **

A Always confirm that the displayed results match the search criteria (e.g., correct cities, dates, and round-trip selection) before proceeding to the next step. 5. If multiple options are displayed, ensure the cheapest or most relevant option is clearly identified and selected as per the task requirements. 6. If a "Back" button fails to function as expected, consider alternative methods to save or exit, such as using a menu or additional options (e.g., "Save as file"). 7. When saving a note as a file, ensure the correct folder and file format are selected before confirming the save. 8. Double-check that the task is fully completed (e.g., the note is saved in the correct location) before marking it as done. 9. If serolling through content does not reveal new information, consider alternative methods to locate the required details, such as using a search or filter function within the app. 10. If the end of a section is reached and the required information is not found, reassess the search criteria or explore other sections of the app for relevant details. 11. When searching for specific items (e.g., dishes, amenities), use keywords or filters to narrow down results and save time. 12. If repetitive actions (e.g., swiping) fall to yield results, pause and evaluate whether the task can be completed using a different approach or if the information is unavailable. 13. When switching between apps, ensure that the context of the task is maintained, and verify that the information gathered in one app aligns with the requirements in the other app. 14. Always confirm the proximity or location details (e.g., using Google Maps) before finalizing a selection, especially when location is a key criterion.

Figure 13: Full list of Tips generated by Mobile-Agent-E (with GPT-40) after self-evolution.