
Mobile-Env: Building Qualified Evaluation Benchmarks for LLM-GUI Interaction

Danyang Zhang¹ Zhennan Shen¹ Rui Xie¹ Situo Zhang¹ Tianbao Xie²
Zihan Zhao¹ Siyuan Chen¹ Lu Chen¹ Hongshen Xu¹ Ruisheng Cao¹ Kai Yu¹

¹Shanghai Jiao Tong University, Shanghai, China

²The University of Hong Kong, Hong Kong, China
zhang-dy20@sjtu.edu.cn

Abstract

The Graphical User Interface (GUI) is pivotal for human interaction with the digital world, enabling efficient device control and the completion of complex tasks. Recent progress in Large Language Models (LLMs) and Vision Language Models (VLMs) offers the chance to create advanced GUI agents. To ensure their effectiveness, there’s a pressing need for qualified benchmarks that provide *trustworthy* and *reproducible* evaluations — a challenge current benchmarks often fail to address. To tackle this issue, we introduce **Mobile-Env**, a comprehensive toolkit tailored for creating GUI benchmarks in the Android mobile environment. Mobile-Env offers an *isolated and controllable setting* for *reliable evaluations*, and accommodates *intermediate instructions and rewards* to reflect real-world usage more naturally. Utilizing Mobile-Env, we collect an open-world task set across various real-world apps and a fixed world set, WikiHow, which captures a significant amount of dynamic online contents for fully controllable and reproducible evaluation. We conduct comprehensive evaluations of LLM agents using these benchmarks. Our findings reveal that even advanced models (e.g., GPT-4V and LLaMA-3) struggle with tasks that are relatively simple for humans. This highlights a crucial gap in current models and underscores the importance of developing more capable foundation models and more effective GUI agent frameworks.¹

1 Introduction

Graphical User Interface (GUI) serves as the main approach for humans to interact with digital devices like mobile phones, personal computers, and web pages. People perform a vast number of tasks through GUIs, such as web browsing, article reading, ticket booking, and file revision. Most of these tasks are time-consuming and repetitive. By employing powerful autonomous agents for these GUI tasks, we can significantly alleviate workloads and boost productivity. Recent Large Language Models (LLMs) and Vision Language Models (VLMs) [16, 14, 15, 1, 4, 26] demonstrate impressive capabilities in understanding both natural and structured texts as well as multimodal contents. They show great potential for interacting with complex multimodal structural GUI environments. The advancement of these foundation models presents a promising opportunity to develop general-purpose GUI agents. Consequently, to fairly assess the capabilities of a wide array of LLM/VLM-based agents, it is essential to build a benchmark for GUI interaction that can thoroughly and reliably evaluate agents in a controllable manner.

Currently, a large portion of benchmarks for GUI interaction are static datasets [21, 3, 13, 11]. These datasets provide fixed ground-truth annotations and evaluate the model’s performance through action

¹Mobile-Env is open-sourced at <https://github.com/X-LANCE/Mobile-Env>. WikiHow task set is publicly available at <https://huggingface.co/datasets/X-LANCE/WikiHow-taskset>.

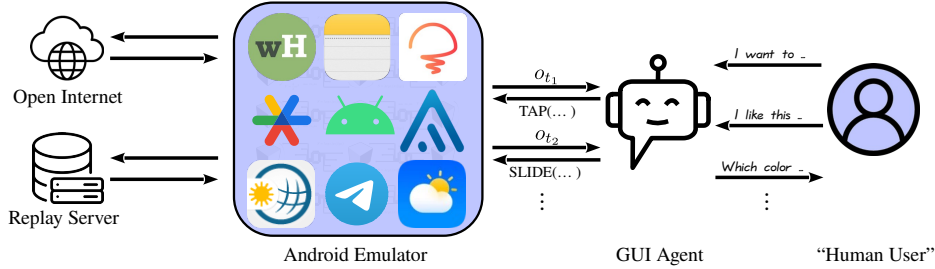


Figure 1: Mobile-Env platform. Mobile-Env comprises an Android emulator and a simple “real user” simulator (the purple sections). The Android emulator can host various real-world apps and run them in either open Internet or strictly controlled environments, *i.e.*, the replay server of crawled resources. The agent interacts with the emulator conditioned by the user’s instructions. Intermediate instructions may be received during execution to supplement task information or user intents.

sequence matching. Therefore, trajectories with alternative routes or detours while still completing the task would be mistakenly marked as failures, leading to wrong evaluations. Interactive benchmarks like AndroidArena [28] adopt similarity of action sequences just like static datasets, and thus fail to assess agent performances adequately. Recent works on LLM-based GUI agents [36, 25, 35] conduct evaluations with the help of LLMs or human judgments. However, their experiments are conducted in completely open environments without any isolation or controls, making the results unreproducible and incomparable.

Furthermore, existing interactive environments lack two crucial features: intermediate rewards and intermediate instructions. Interactive environments should support more agent designing methods, such as Reinforcement Learning (RL) [6]. In these cases, intermediate rewards are essential for agents to learn efficiently in GUI environments with complex observations and extensive action space. Additionally, it is unrealistic for users to provide a comprehensive command containing all task details at once. Therefore, intermediate instructions should be provided throughout the task based on information from newly encountered screens, creating more natural and practical use cases.

In general, a qualified evaluation benchmark for GUI interaction is expected to have several features: 1) **reliable evaluation** that assessing the execution result rather than the execution trajectory as multiple feasible trajectories may exist, 2) **isolated and controllable environments** to prevent agents’ latent dangerous actions from corrupting real-world environments and guarantee reproducible evaluation and fair comparison, 3) **intermediate rewards** for better supporting interactive learning methods, and 4) **intermediate instructions** to face more natural and realistic use cases.

Towards building qualified evaluation benchmarks for GUI interaction, we developed the Mobile-Env platform as a comprehensive toolkit for GUI benchmark construction. Noticed that mobiles hold the most market share among consumer electronic devices², we build Mobile-Env based on the Android™ mobile platform. Mobile-Env offers an isolated interactive environment and designs a flexible framework to reliably evaluate the widest range of tasks. Additionally, intermediate rewards in Mobile-Env enable efficient reinforcement learning, and the intermediate instructions facilitate the development of conversational agents. With the help of Mobile-Env, qualified evaluation benchmarks can be built much more easily.

Based on Mobile-Env, we design 74 tasks across a number of real-world apps. We further build a task set, WikiHow, containing 150 tasks generated from 13 templates, by fixing a remarkable scale of dynamic online contents, achieving absolutely controllable and reproducible evaluation. Extensive experiments are conducted on these designed tasks with current advanced LLMs and VLMs. It is found that even the most advanced models struggle with tasks that are relatively simple for humans. These results reveal the deficiency of current LLMs and VLMs as reliable and effective GUI assistants, underscoring the necessity for developing more capable foundation models and knowledge-augmented agents.

²<https://gs.statcounter.com/platform-market-share/desktop-mobile-tablet/worldwide/2023>

Table 1: Comparison of Mobile-Env with existing GUI benchmarks. The table shows the support of the benchmarks/benchmark toolkits for result-based evaluation, truly controllable environments, intermediate rewards, and intermediate instructions.

	Benchmark (Toolkit)	Result-based Evaluation	Controllable Environments	Intermediate Rewards	Intermediate Instructions
Static Dataset	PixelHelp [8]	✗	-	-	✗
	META-GUI [21]	✗	-	-	✓
	Mind2Web [3]	✗	-	-	✗
	AitW [17]	✗	-	-	✗
	WebLINX [11]	✗	-	-	✓
Interactive Environment	MiniWoB++ [10]	✓	✓	✗	✗
	AndroidEnv [23]	✓	✗	✓	✗
	WebShop [31]	✓	✓	✗	✗
	WebArena [37]	✓	✓	✗	✗
	VisualWebArena [7]	✓	✓	✗	✗
	AndroidArena [28]	✗	✗	✗	✗
	OSWorld [27]	✓	✗	✗	✗
	AndroidWorld [18]	✓	✗	✗	✗
	Mobile-Env (Ours)	✓	✓	✓	✓

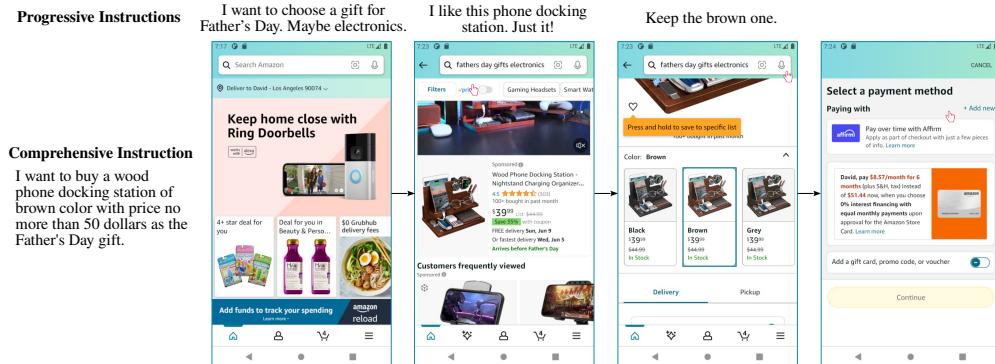


Figure 2: Progressive instructions vs. comprehensive instruction. A single comprehensive instruction that covers every detail is usually verbose and unnatural. In real use cases, the user's intent may not be fully clear at the beginning but becomes more defined during the execution process.

2 Background

2.1 Problem formulation

As illustrated in Fig. 1, a GUI agent A receives the human user's instructions and completes the task in a GUI system (*i.e.*, the Android emulator in the figure). The "human user" and the GUI system are regarded as the environment holding hidden intents/internal states s_t , which the agent should deduce via observations. The observation o_t the agent receives at each step consists of the screen information from the GUI system and an optional intermediate instruction from the "human user". Then it predicts an action a_t to operate the GUI system or to respond to the "human user", causing a state transition (s_t, a_t, s_{t+1}) . An optional reward can be given to assess the transition and can be used to improve the agent. Observation receiving and action taking iterate alternatively in this manner until the task comes to success or failure.

2.2 Qualified benchmarks

This subsection will discuss the status quo of existing GUI benchmarks and the proposed features of a qualified benchmark. An overall comparison is available in Tab. 1.

Reliable evaluation Existing GUI interaction benchmarks put a lot of efforts into evaluating the agent’s executions. Most static benchmarks [8, 21, 3, 17, 13, 11] simply perform step-wise evaluation by comparing the predicted actions with the reference action trajectory. However, there are usually multiple feasible solutions for a specific task goal. The agent may also conduct explorations not directly leading to task completion. In order to conduct reliable evaluations, the benchmark should inspect the execution *result* rather than the execution *process* [37, 7]. Specifically, the task goal can always be viewed as a subset of hidden states. As long as the agent reaches these states, its execution will be deemed successful. The hidden final state is supposed to be estimated through a series of system signals. Usually, a single type of signal cannot cover all the state transitions, *e.g.*, several changes like form submission can only be reflected on GUI, but not in the system logs adopted in Toyama et al. [23] or files and databases adopted in Rawles et al. [18]. Therefore, diverse system signals should be leveraged to comprehensively assess the final states [27]. Mobile-Env exploits a wide range of signal types, including screen text, View Hierarchy (VH)³, system logs, and agent’s Response to Human User (RHU). Furthermore, Mobile-Env designs flexible mechanisms to combine multiple signals, achieving reliable and flexible evaluation.

Isolated & controllable environments Due to the lack of qualified benchmarks, recent work on LLM-based GUI agents [5, 36, 25, 35] conducts evaluations with real devices in fully open environments without any isolation or control. This approach leads to several issues: safety risks, unreproducible evaluations, and unfair comparisons. *Safety risks* are referred to as the fact that without environment **isolation**, the agent may mistakenly corrupt real-world devices or data resources. On the other hand, many apps rely on online contents that vary at different time and locations. Therefore, environments built on such apps will result in *unreproducible evaluation*. Meanwhile, these factors will also affect the actual difficulty of the tasks in the environment and result in *unfair comparisons* between agents. For instance, an agent might easily find a Panda Express⁴ in the USA but not in China. Consequently, the only way to set up reproducible and fair evaluation is to build a truly **controllable** environment by fixing the dynamic online contents and replaying during evaluations [19, 10, 37]. To enable common data replay for mobile apps, we reverse-engineered and developed a suite of solutions for replay server usage and app certificate unpinning.

Intermediate rewards Reinforcement Learning (RL) is a method of great potential for building intelligent agents. Extensive work is trying to leverage RL to build digital agents [6, 30, 38, 34]. However, a common challenge in RL is reward sparsity, which is exacerbated by the increasingly complex observations and extensive action space in GUI environments. One way to alleviate this problem is to introduce intermediate rewards [20, 2, 12]. Mobile-Env designs a listener-style evaluator and produces intermediate rewards safely for GUI interaction tasks.

Intermediate instructions As shown in Fig. 2, a comprehensive task instruction covering all details is usually lengthy and verbose. It is impractical for real users to consider all variable factors at the outset. Realistic task instructions should be progressive. Initially, there is a simple instruction that outlines the user’s basic intent. As the task progresses and new information emerges, the user’s intent becomes clearer, and more details are provided. Moreover, supporting intermediate instructions is crucial for evaluating conversational GUI agents [21, 11], enabling agents and users to communicate about task details and progress. To our knowledge, Mobile-Env is the first interactive environment that supports intermediate instructions, allowing for the evaluation of conversational GUI agents.

3 Mobile-Env platform

Mobile-Env offers a uniform framework to host various GUI environments (§ 3.1) and build qualified GUI benchmarks (§ 3.2). As it is a common way to “work with the emulators thus building these environments”, we will introduce the internal components and implementation details in the supplementary materials, and leave this section focusing on the use features of Mobile-Env.

³A tree-shape representation of the Android screens, similar to HTML of the web pages.

⁴A popular Chinese food restaurant in the USA.

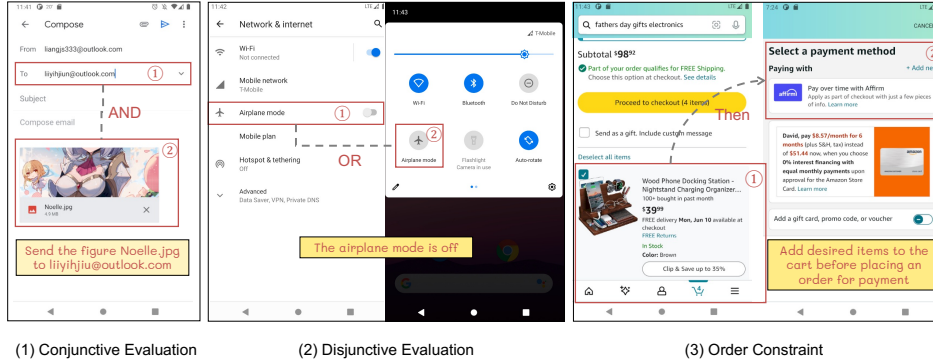


Figure 3: Mobile-Env utilizes three methods (*i.e.*, conjunctive evaluation, disjunctive evaluation, and order constraints) for combining signals to achieve flexible evaluation.

3.1 Universal hosting of GUI environments

Mobile-Env is capable of hosting most real-world Android apps and offers uniform interfaces for all kinds of apps and mobile interaction tasks. To be specific, the observation of Mobile-Env environments include both screenshots and VH (View Hierarchy) XML. The raw GUI action space of Mobile-Env is pixel-wise atomic TOUCH (finger down), LIFT (finger up), and TEXT (token input), ensuring a universal and fine-grained interface for interacting with various apps. Additionally, we provide a convenient way to use environment wrappers to modify the action space; for example, appending a LIFT to a series of TOUCH actions to create a TAP. The RHU (Response to Human User) action is also available in Mobile-Env. As shown in Fig. 1, the agent can interact with Mobile-Env environments either in an open-world setting or within a strictly controlled world conditioned by multi-turn instructions. This makes Mobile-Env a universal platform for various mobile interactive environments and agents.

3.2 Creation of qualified GUI benchmarks

Mobile-Env is used to build qualified GUI interaction benchmarks. Its first notable feature lies in the configuration setup: all the setup procedures, evaluation methods, and reward mechanisms are configured through an external task configuration file. This file adopts the text format of Protocol Buffer⁵, which is easier to comprehend and modify compared to complicated and overwrapped programming codes. The subsequent part will introduce Mobile-Env’s methodology towards building qualified benchmarks by (1) implementing isolated and controllable environments, and (2) ensuring accurate task state estimation for reliable evaluation and immediate rewards & instructions.

3.2.1 Setup of isolated & controllable environments

Mobile-Env sets up an environment using a predefined Android OS (Operating System) image and executes a sequence of ADB (Android Debug Bridge) commands defined in the task configuration file. This process ensures that the environment is properly prepared for the task before awaiting the agent’s execution. If a fully controlled “online” environment needs to be established, a replay server for the crawled app data should be launched in advance. However, direct replay will face the challenge of certificate pinning in real-world apps, which prevents them from trusting SSL certificates issued by the replay proxy and leads to data replay failure. To handle this problem, we implemented three solutions according to the conclusions of reverse engineering, and verified these solutions on a series of common real-world apps. Further details on these solutions are provided in the supplementary.

3.2.2 Accurate task state estimation

Accurate task state estimation is essential for both reliable evaluation and to guide the agent with intermediate rewards and instructions at the right time. Mobile-Env ensures this accuracy by utilizing a variety of signals including screen text, VH, system log, and RHU. Moreover, it introduces three

⁵<https://protobuf.dev/reference/protobuf/textformat-spec/>

Table 2: Task instruction examples of the open world task set

Instruction	Involved Apps	IR*	II†
Turn on airplane mode for my phone.	System Settings	✗	✗
Check my TOTP‡ for GitHub and tell me the code.	Google Authenticator	✗	✗
Help me to order a cup of my favorite milk tea of Coco (the store name).	Ele.me	✓	✓
I have several funds and bills to be noted in OpenMoneyBox. Now open the app, enter the wizard, and help me to note the bills.	OpenMoneyBox	✓	✓
Turn off my alarm in the morning if it will rain tomorrow, otherwise leave it on.	Weather, Alarm Clock	✗	✗

* Intermediate Rewards † Intermediate Instructions ‡ Time-based One-Time Password

strategies for combining these signals to enhance the flexibility of task state estimation (see Fig. 3): (1) **Conjunctive evaluation**: multiple signals must occur simultaneously, e.g., an e-mail is ready to be sent only if the recipient and contents are properly filled in. (2) **Disjunctive evaluation**: multiple alternative signals are combined to avoid missing evidence of success. For instance, airplane mode on a mobile can be toggled from either the Quick Settings menu or the System Settings app. The final pages differ but both indicate whether airplane mode is on or off. (3) **Order constraint**: complex tasks consist of multiple stages, with some serving as prerequisites for others. For example, desired items must be added to the cart before placing an order for payment. By employing the aforementioned methods, Mobile-Env can achieve more accurate state estimation. Additional examples of task evaluation are detailed in the supplementary materials.

Furthermore, Mobile-Env adopts a “listener-style” approach for capturing task signals without disrupting the runtime state or the agent’s actions. This design is different from the “post-processing” evaluators as described by Xie et al. [27], which might involve actions like saving an open file or closing an app that could potentially alter the environment state. The listener style adopted by Mobile-Env guarantees safe and correct state evaluations, as well as enables intermediate rewards and instructions, without interfering with the agent’s ongoing execution.

4 Task sets

With the help of Mobile-Env, we collected two sets of GUI interaction tasks, one in the open world, and another in a fixed controlled world. This section will briefly introduce these two task sets.

4.1 Open world set

We define 74 GUI interaction tasks using real-world apps running directly in an open Internet environment to evaluate the agents’ raw performance. Following AitW [17] and AndroidWorld [18], we include app installation/uninstallation, clock settings, and Android system settings in the test set. Additionally, we select four popular apps and define a series of tasks based on their functions. The chosen apps are the instructional app WikiHow, the food delivery app Ele.me, and the review platforms Douban and IMDB.⁶ We also define several tasks involving multiple apps to test the agents’ ability to navigate across different applications. Examples of task instructions are shown in Tab. 2, and the full list of task instructions can be found in the supplementary materials.

4.2 Fixed world set (WikiHow)

Aiming at providing truly comparable results, we further collect a fixed world set based on the WikiHow app. WikiHow is an app hosting millions of “how-to” articles about step-by-step everyday tips. We crawled 107,448 pages containing 856,045 text and multimedia resources from the WikiHow website to establish the fixed controlled environment. Then, we craft 13 task configuration templates and generate 150 tasks for testing. The template-based instructions are further rewritten with the help

⁶WikiHow: <https://www.wikihow.com/Main-Page>. Ele.me: <https://www.ele.me/>. Douban: <https://www.douban.com/>. IMDB: <https://www.imdb.com/>.

of ChatGPT⁷ to increase the diversity of expression. Multiple rewritten candidates are recorded for each instruction and will be replaced randomly at runtime. An example of sequential instructions from the fixed WikiHow task set is shown in Tab. 3. This task includes a series of sequential instructions guiding the agent to browse app contents. An intermediate reward is given upon completing each sub-instruction, and the episode is considered successful only when the final instruction is completed. Details of app data crawling, template crafting and instantiation, and instruction replacing are provided in the supplementary materials.

To draw deeper insights into the agents’ capabilities, we split the fixed world set into three categories: 1) 59 *cross-page tasks* requiring pure navigation among pages, 2) 51 *in-page tasks* engaging in several in-page operations, like bookmarking or voting for an article, and 3) 40 *QA tasks* requesting the agent to give a summarized answer according to the article contents, *e.g.*, “read the article and summarize the needed items for the task”.

Table 3: An example of sequential instructions in the fixed WikiHow task set

- | |
|--|
| 1. How can I hide my ear gauges? |
| 2. <i>How to hide gauges</i> may be helpful. |
| 3. Share the article to other apps. |

5 Experiments

5.1 Experiment settings

We experiment with current leading LLMs and VLMs across three categories: (1) *Closed-source LLMs* like GPT-3.5, GPT-4, and Claude-3, (2) *Open-source LLMs* including LLaMA-2, LLaMA-3, and AgentLM of varying sizes, and (3) *VLMs* such as GPT-4V and Claude-3-Opus. Text-based LLMs take a simplified HTML derived from the VH observation as input. For VLMs, we provide the screenshot and Set-of-Marks (SoM) [29]. Set-of-Marks are indexed visual tags annotating the essential GUI elements on the image, helping the VLM to ground in the screenshot. We prompt the agents to first perform reasoning before the actual action, following Yao et al. [32]. The average reward (Rwd) through episode and success rate (SR, in %) are reported as the metrics.

Table 4: Results on open world set

Model	Rwd	SR
Claude-3-Opus	0.86	32.43
Claude-3-Opus (V)	0.25	2.70
GPT-4	1.09	43.24
GPT-4V	0.18	3.04
AgentLM-70B	0.77	13.51
LLaMA-3-70B	1.11	40.54

5.2 Results on open world set

The results on the open world set are shown in Tab. 4, demonstrating the low performances across all the models, especially the VLMs (SRs of 2.70% for Claude-3-Opus and 3.04% for GPT-4V). This implies the deficiency of current LLMs and VLMs as effective GUI assistants. Nevertheless, it is surprising to see a comparable performance of LLaMA-3-70B (SR of 40.54%) with text GPT-4 (SR of 43.24%), showing the incredible improvement of open-source models. Through carefully looking into the trajectories, we attribute the deficiency of the models, especially the VLMs, to the lack of the essential capabilities like common-sense GUI understanding and output format following.

Common-sense knowledge of GUI understanding GUI design follows user-friendly principles, employing familiar visual elements such as icons, buttons, and menus, which users can quickly recognize and interpret. In contrast, our experiments show that current LLMs and VLMs often lack this common sense and struggle to understand GUIs, especially on information-intensive screens like weather apps displaying verbose forecasts for many future days. For example, few agents can recognize the contents shown in File Manager and share a file from it. And GPT-4 doesn’t know how to display the receiving time of an e-mail. Specifically, we find that the VLMs perform much worse than text LLMs, as the screen texts and text prompts of icons are not as explicitly available in the screenshot as in the HTML. Thus, the VLMs face more difficulties in understanding the texts, icons, and functions of the current screen, leading to terribly low performance.

⁷<https://chat.openai.com/>

Table 5: Experiment results of text LLMs on WikiHow task set

Model	Rwd	SR	Cross-Page		In-Page		QA	
			Rwd	SR	Rwd	SR	Rwd	SR
Claude-3-Opus	1.94	40.67	1.61	37.29	2.51	54.90	1.71	27.50
GPT-3.5	1.79	37.33	1.49	38.98	2.29	49.02	1.58	20.00
GPT-3.5 (1-turn)	1.81	28.67	1.67	40.68	2.16	27.45	1.58	12.50
GPT-3.5-instruct	1.76	36.67	1.71	44.07	2.02	35.29	1.50	27.50
GPT-4	2.01	43.33	1.95	61.02	2.35	50.98	1.65	7.50
LLaMA-2-7B	0.65	4.00	0.81	10.17	0.78	0.00	0.25	0.00
LLaMA-2-13B	1.08	11.33	1.14	27.12	1.37	0.00	0.63	2.50
LLaMA-2-70B	1.49	20.00	1.44	30.51	1.65	7.84	1.34	20.00
LLaMA-3-8B	1.74	30.67	1.58	35.59	1.90	23.53	1.77	32.50
LLaMA-3-70B	2.10	47.33	1.83	49.15	2.47	54.90	2.02	35.00

Format following capabilities Language model based agents are expected to output well-formatted texts so that the actions can be easily parsed and executed. Our experiments show that even the most capable LLMs and VLMs struggle to adhere to the output format strictly. Over half the predictions of GPT-4 and Claude-3-Opus contain redundant “new line” characters. When input with screenshots, ~26.40% outputs of GPT-4V and ~10.87% outputs of Claude-3-Opus comprise only thoughts without actual actions. This problem has become a significant bottleneck of agent’s performance.

5.3 Results of text LLMs on fixed world set

As a supplement to results on the open world set, we further conduct experiments on the fixed world set to provide truly reproducible results. The results of text LLMs on the WikiHow task set are presented in Tab. 5. It is observed that even the best performance achieves a success rate of only 47.33%, which highlights the deficiency of current LLMs on GUI interactions again.

Performances across different task types Overall, the QA tasks seem more difficult than the other types of tasks. This is not surprising, as the QA tasks require the agent to scroll up and down to read the whole article and extract the useful text contents out of the HTML markups. An example of the HTML representation on the article page is shown in Fig. 4. Thus, this task is harder than simply seeking a specific button or hyperlink. Meanwhile, it is noticed that many models perform worse on the in-page tasks than on the cross-page tasks. By inspecting the interaction trajectories, this is attributed to that several in-page tasks necessitate actions like long-range scrolling or checking into dropdown menus. Successfully executing these actions typically requires a strong belief in their effectiveness. However, many agents, lacking such confidence, usually try to click something else during scrolling, aiming at reaching the target more quickly. Besides, there may be shortcuts that could simplify task completion, such as clicking a reference list link to quickly navigate to the bottom of a page instead of long-term scrolling. However, few models possess such *a priori* knowledge and manage to complete the tasks in this manner. How to establish appropriate belief and inject/acquire *a priori* knowledge of the underlying environment may be valuable directions to explore for automatic agent design. These insights are further discussed with the case shown in the supplementary.

Multi-turn prompting vs. single-turn prompting We conducted an experiment to investigate the impact of the prompting format for chat models. Multi-turn prompting is referred to as splitting the input exemplars into observation and action and assigning observation the “user” role and action the “assistant” role. In contrast, single-turn prompting is referred to as assigning all the input prompts (except the “system” part) the role “user” and expecting a 1-turn response from the LLM

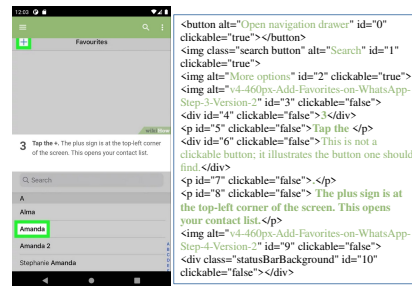


Figure 4: HTML input of QA tasks. The article cannot be perceived as whole and the text contents are interleaved with HTML markups.

assistant. As depicted by the results of GPT-3.5 in Tab. 5, the single-turn format dramatically degrades the performance. We speculate that multi-turn prompting combining the input-output logic with user-assistance distinction makes it clearer for the model to learn the relation between inputs and outputs.

5.4 Results of VLMs on fixed world set

We evaluated the vision-based agents on 40 tasks selected from the fixed world task set. As shown in Tab. 6, even the strongest VLM, GPT-4V, achieved only a 10.00% success rate. Except for the reasons discussed in § 5.2, it is additionally attributed to that current VLMs struggle to accurately ground on GUI screenshots. Notably, the SoM approach mitigates this challenge and improves VLM’s performance by asking the VLMs to predict the index of a tagged screen element instead of its precise coordinates. The result demonstrates there is significant room for improvement in current VLMs’ ability to reason and correctly ground on the GUI screenshots.

Table 6: Experiment results of VLM-based agents on WikiHow task set. The results are obtained with/(out) SoM and a different number of exemplars (Exem). “Claude” denotes Claude-3-Opus.

Model	SoM?	#Exems	Rwd	SR
Claude	✓	2	0.28	0.00
GPT-4V	✓	2	0.90	10.00
GPT-4V	✗	2	0.33	3.00
GPT-4V	✓	0	0.70	5.00

6 Conclusion

Aiming at filling in the blank of a qualified GUI interaction benchmark, this work builds Mobile-Env as a comprehensive toolkit to simplify the construction of qualified benchmarks. Mobile-Env offers an isolated interactive environment and designs a flexible framework to produce reliable intermediate rewards and instructions. With the help of Mobile-Env, two task sets are collected to draw insights into the performances of LLM-based agents and provide truly reproducible and comparable results. The experiments reveal the deficiency of current models in serving as GUI agents, underscoring the necessity to build more capable foundation models and more effective and robust agent frameworks augmented with knowledge of GUI understanding and operation.

References

- [1] AI@Meta. Llama 3 model card. 2024. URL https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md.
- [2] Jake Bruce, Ankit Anand, Bogdan Mazoure, and Rob Fergus. Learning about progress from experts. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL <https://openreview.net/pdf?id=sKc6fgc1zs>.
- [3] 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. *CoRR*, abs/2306.06070, 2023. doi: 10.48550/ARXIV.2306.06070. URL <https://doi.org/10.48550/arXiv.2306.06070>.
- [4] Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. GLM: general language model pretraining with autoregressive blank infilling. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 320–335. Association for Computational Linguistics, 2022. doi: 10.18653/v1/2022.acl-long.26. URL <https://doi.org/10.18653/v1/2022.acl-long.26>.
- [5] Difei Gao, Lei Ji, Zechen Bai, Mingyu Ouyang, Peiran Li, Dongxing Mao, Qinchen Wu, Weichen Zhang, Peiyi Wang, Xiangwu Guo, Hengxu Wang, Luwei Zhou, and Mike Zheng Shou. ASSISTGUI: task-oriented desktop graphical user interface automation. *CoRR*, abs/2312.13108, 2023. doi: 10.48550/ARXIV.2312.13108. URL <https://doi.org/10.48550/arXiv.2312.13108>.

- [6] Peter C. Humphreys, David Raposo, Tobias Pohlen, Gregory Thornton, Rachita Chhaparia, Alistair Muldal, Josh Abramson, Petko Georgiev, Adam Santoro, and Timothy P. Lillicrap. A data-driven approach for learning to control computers. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato, editors, *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 9466–9482. PMLR, 2022. URL <https://proceedings.mlr.press/v162/humphreys22a.html>.
- [7] Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Chong Lim, Po-Yu Huang, Graham Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and Daniel Fried. Visualwebarena: Evaluating multimodal agents on realistic visual web tasks. *CoRR*, abs/2401.13649, 2024. doi: 10.48550/ARXIV.2401.13649. URL <https://doi.org/10.48550/arXiv.2401.13649>.
- [8] Yang Li, Jiacong He, Xin Zhou, Yuan Zhang, and Jason Baldridge. Mapping natural language instructions to mobile UI action sequences. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel R. Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 8198–8210. Association for Computational Linguistics, 2020. doi: 10.18653/v1/2020.acl-main.729. URL <https://doi.org/10.18653/v1/2020.acl-main.729>.
- [9] Jimmy Lin, Xueguang Ma, Sheng-Chieh Lin, Jheng-Hong Yang, Ronak Pradeep, and Rodrigo Nogueira. Pyserini: A python toolkit for reproducible information retrieval research with sparse and dense representations. In Fernando Diaz, Chirag Shah, Torsten Suel, Pablo Castells, Rosie Jones, and Tetsuya Sakai, editors, *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021*, pages 2356–2362. ACM, 2021. doi: 10.1145/3404835.3463238. URL <https://doi.org/10.1145/3404835.3463238>.
- [10] Evan Zheran Liu, Kelvin Guu, Panupong Pasupat, Tianlin Shi, and Percy Liang. Reinforcement learning on web interfaces using workflow-guided exploration. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=ryTp3f-0->.
- [11] Xing Han Lu, Zdenek Kasner, and Siva Reddy. Weblinx: Real-world website navigation with multi-turn dialogue. *CoRR*, abs/2402.05930, 2024. doi: 10.48550/ARXIV.2402.05930. URL <https://doi.org/10.48550/arXiv.2402.05930>.
- [12] Chang Ma, Junlei Zhang, Zhihao Zhu, Cheng Yang, Yujiu Yang, Yaohui Jin, Zhenzhong Lan, Lingpeng Kong, and Junxian He. Agentboard: An analytical evaluation board of multi-turn LLM agents. *CoRR*, abs/2401.13178, 2024. doi: 10.48550/ARXIV.2401.13178. URL <https://doi.org/10.48550/arXiv.2401.13178>.
- [13] Runliang Niu, Jindong Li, Shiqi Wang, Yali Fu, Xiyu Hu, Xueyuan Leng, He Kong, Yi Chang, and Qi Wang. Screenagent: A vision language model-driven computer control agent. *CoRR*, abs/2402.07945, 2024. doi: 10.48550/ARXIV.2402.07945. URL <https://doi.org/10.48550/arXiv.2402.07945>.
- [14] OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. doi: 10.48550/arXiv.2303.08774. URL <https://doi.org/10.48550/arXiv.2303.08774>.
- [15] OpenAI. Gpt-4v(ision) system card, 2023. URL https://cdn.openai.com/papers/GPTV_System_Card.pdf. Accessed at Dec 16, 2023.
- [16] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *NeurIPS*, 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html.

- [17] Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy P. Lillicrap. Android in the wild: A large-scale dataset for android device control. *CoRR*, abs/2307.10088, 2023. doi: 10.48550/ARXIV.2307.10088. URL <https://doi.org/10.48550/arXiv.2307.10088>.
- [18] Christopher Rawles, Sarah Clinckemaulle, 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.
- [19] Tianlin Shi, Andrej Karpathy, Linxi Fan, Jonathan Hernandez, and Percy Liang. World of bits: An open-domain platform for web-based agents. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 3135–3144. PMLR, 2017. URL <http://proceedings.mlr.press/v70/shi17a.html>.
- [20] Maayan Shvo, Zhiming Hu, Rodrigo Toro Icarte, Iqbal Mohamed, Allan D. Jepson, and Sheila A. McIlraith. Appbuddy: Learning to accomplish tasks in mobile apps via reinforcement learning. In Luiza Antonie and Pooya Moradian Zadeh, editors, *Proceedings of the 34th Canadian Conference on Artificial Intelligence, Canadian AI 2021, online, May 2021*. Canadian Artificial Intelligence Association, 2021. URL <https://caiac.pubpub.org/pub/5agsashz>.
- [21] Liangtai Sun, Xingyu Chen, Lu Chen, Tianle Dai, Zichen Zhu, and Kai Yu. META-GUI: towards multi-modal conversational agents on mobile GUI. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pages 6699–6712. Association for Computational Linguistics, 2022. doi: 10.18653/V1/2022.EMNLP-MAIN.449. URL <https://doi.org/10.18653/v1/2022.emnlp-main.449>.
- [22] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. *CoRR*, abs/2307.09288, 2023. doi: 10.48550/ARXIV.2307.09288. URL <https://doi.org/10.48550/arXiv.2307.09288>.
- [23] Daniel Toyama, Philippe Hamel, Anita Gergely, Gheorghe Comanici, Amelia Glaese, Zafarali Ahmed, Tyler Jackson, Shibl Mourad, and Doina Precup. Androidenv: A reinforcement learning platform for android. *CoRR*, abs/2105.13231, 2021. URL <https://arxiv.org/abs/2105.13231>.
- [24] Bryan Wang, Gang Li, and Yang Li. Enabling conversational interaction with mobile UI using large language models. *CoRR*, abs/2209.08655, 2022. doi: 10.48550/ARXIV.2209.08655. URL <https://doi.org/10.48550/arXiv.2209.08655>.
- [25] 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. *CoRR*, abs/2401.16158, 2024. doi: 10.48550/ARXIV.2401.16158. URL <https://doi.org/10.48550/arXiv.2401.16158>.
- [26] Weihang Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, Jiazheng Xu, Bin Xu, Juanzi Li, Yuxiao Dong, Ming Ding, and Jie Tang. Cogvlm: Visual expert for pretrained language models. *CoRR*, abs/2311.03079, 2023. doi: 10.48550/ARXIV.2311.03079. URL <https://doi.org/10.48550/arXiv.2311.03079>.

- [27] 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. *CoRR*, abs/2404.07972, 2024. doi: 10.48550/ARXIV.2404.07972. URL <https://doi.org/10.48550/arXiv.2404.07972>.
- [28] Mingzhe Xing, Rongkai Zhang, Hui Xue, Qi Chen, Fan Yang, and Zhen Xiao. Understanding the weakness of large language model agents within a complex android environment. *CoRR*, abs/2402.06596, 2024. doi: 10.48550/ARXIV.2402.06596. URL <https://doi.org/10.48550/arXiv.2402.06596>.
- [29] Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chunyuan Li, and Jianfeng Gao. Set-of-mark prompting unleashes extraordinary visual grounding in GPT-4V. *CoRR*, abs/2310.11441, 2023. doi: 10.48550/ARXIV.2310.11441. URL <https://doi.org/10.48550/arXiv.2310.11441>.
- [30] Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents. In *NeurIPS*, 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/82ad13ec01f9fe44c01cb91814fd7b8c-Abstract-Conference.html.
- [31] Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents. *CoRR*, abs/2207.01206, 2022. doi: 10.48550/arXiv.2207.01206. URL <https://doi.org/10.48550/arXiv.2207.01206>.
- [32] 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/pdf?id=WE_vluYUL-X.
- [33] Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao Liu, Yuxiao Dong, and Jie Tang. Agenttuning: Enabling generalized agent abilities for llms. *CoRR*, abs/2310.12823, 2023. doi: 10.48550/ARXIV.2310.12823. URL <https://doi.org/10.48550/arXiv.2310.12823>.
- [34] Yuexiang Zhai, Hao Bai, Zipeng Lin, Jiayi Pan, Shengbang Tong, Yifei Zhou, Alane Suhr, Saining Xie, Yann LeCun, Yi Ma, et al. Fine-tuning large vision-language models as decision-making agents via reinforcement learning. *arXiv preprint arXiv:2405.10292*, 2024.
- [35] 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. *CoRR*, abs/2402.07939, 2024. doi: 10.48550/ARXIV.2402.07939. URL <https://doi.org/10.48550/arXiv.2402.07939>.
- [36] Chi Zhang, Zhao Yang, Jiaxuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. Appagent: Multimodal agents as smartphone users. *CoRR*, abs/2312.13771, 2023. doi: 10.48550/ARXIV.2312.13771. URL <https://doi.org/10.48550/arXiv.2312.13771>.
- [37] Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. Webarena: A realistic web environment for building autonomous agents. *CoRR*, abs/2307.13854, 2023. doi: 10.48550/ARXIV.2307.13854. URL <https://doi.org/10.48550/arXiv.2307.13854>.
- [38] Yifei Zhou, Andrea Zanette, Jiayi Pan, Sergey Levine, and Aviral Kumar. Archer: Training language model agents via hierarchical multi-turn RL. *CoRR*, abs/2402.19446, 2024. doi: 10.48550/ARXIV.2402.19446. URL <https://doi.org/10.48550/arXiv.2402.19446>.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#)

- (b) Did you describe the limitations of your work? [Yes] The limitations and future work are discussed in the supplementary.
 - (c) Did you discuss any potential negative societal impacts of your work? [Yes] The potential societal impacts are discussed in the supplementary.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
- (a) Did you state the full set of assumptions of all theoretical results? [N/A]
 - (b) Did you include complete proofs of all theoretical results? [N/A]
3. If you ran experiments (e.g. for benchmarks)...
- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] The Mobile-Env code repository, WikiHow task set, and the crawled data of the WikiHow website are open-sourced at GitHub and Hugging Face. The links are attached to the main paper. The full prompts to reproduce the experiment results are represented in the supplementary as well.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [N/A]
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [No] Duplicated experiments are not conducted owing to the limit of budget.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See the supplementary.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- (a) If your work uses existing assets, did you cite the creators? [Yes]
 - (b) Did you mention the license of the assets? [Yes] The license is mentioned in the supplementary.
 - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] The links to the Mobile-Env code repository and WikiHow task set are attached in the main paper.
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

A Distribution & License

The proposed Mobile-Env platform and WikiHow task set are open-sourced under Apache-2.0 license. Mobile-Env is released on <https://github.com/X-LANCE/Mobile-Env> and WikiHow task set is distributed through <https://huggingface.co/datasets/X-LANCE/WikiHow-taskset>. Both Mobile-Env and the task sets will be maintained continuously. The experiment codes are open-sourced at <https://github.com/OpenDFM/mobile-env-expe>.

A non-exhaustive list of artifacts used in development of Mobile-Env and WikiHow task set includes: AndroidEnv [23], mitmproxy⁸, Pyserini [9], WikiHow⁹, *etc.* These artifacts are released under

⁸<https://mitmproxy.org>

⁹<https://www.wikihow.com>

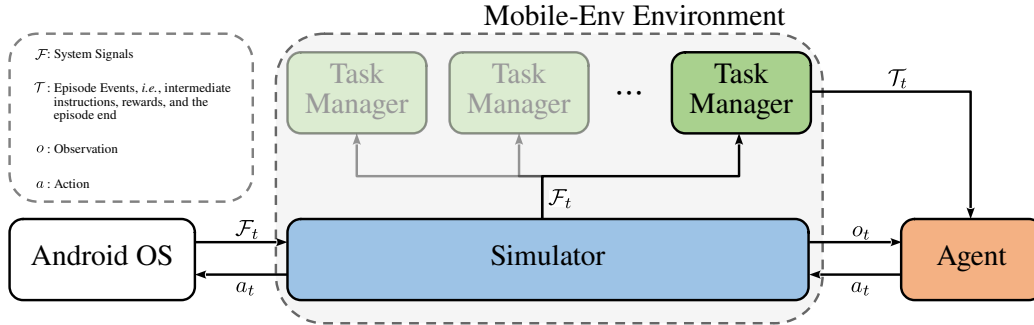


Figure 5: Main components of Mobile-Env. The simulator serves as a proxy to the Android OS, offers observations and accepts actions from the agent. The task manager estimates the task state through various types of system signals and produce episode events. Multiple task managers can be combined with the same simulator to reduce overhead, and only one can be active for interaction at one time.

licenses Apache-2.0, MIT, CC-BY-NC-SA 3.0, *etc.* A non-exhaustive list of artifacts used in the experiments includes: LLaMA-2 [22], AgentLM [33], LLaMA-3 [1], vLLM¹⁰, Text-Generation-Inference¹¹, NLTK¹², EasyOCR¹³, vision-ui¹⁴, *etc.* They are released under licenses LLaMA 2, LLaMA 3, HFOILv1, MIT, Apache-2.0, *etc.* The authors claim that the usage completely obey the licenses.

B Implementation details of Mobile-Env

Mobile-Env is designed as a universal interaction platform hosting various apps and different tasks to facilitate qualified evaluation of LLM/VLM-based or conventional GUI agents. A task is corresponding to a specific goal and several configurations, which are defined by a task configuration file. Based on Mobile-Env, environment and task designers can easily build qualified environments and tasks. Then, agent designers can train and evaluate their GUI agents with various environments and tasks held on Mobile-Env.

As shown in Fig. 5, a runtime instance of Mobile-Env consists of two types of components:

Simulator The simulator wraps Android Emulator^{TM15} and provides the agent with a suite of interaction interfaces with the AndroidTM OS. The simulator predefines basic observation space and action space.

Task manager A task manager is corresponding to a specific task goal. During the launch of Mobile-Env instance, the task manager will be loaded by reading and parsing a task configuration file. Then, during the episode, the task manager will produce intermediate instructions and rewards and indicate the end of episode (either failure or success) at the correct time according to the loaded task configuration.

At runtime, multiple task managers can be combined with one simulator to prevent the overhead of multiple simulator instances. Note that there will be only one task manager active for the current episode at one time.

During interaction, the agent will receive the observation from the simulator and the task events (*i.e.*, intermediate instructions, rewards, and the episode end) from the active task manager. The action decision will then be made accordingly and sent to the OS through the simulator.

¹⁰<https://docs.vllm.ai/en/latest/index.html>

¹¹<https://github.com/huggingface/text-generation-inference>

¹²<https://www.nltk.org/>

¹³<https://www.jaided.ai/easyocr/>

¹⁴<https://github.com/Meituan-Dianping/vision-ui>

¹⁵<https://developer.android.com/about>

Table 7: Time delay of acquiring different types of observation and performing different actions

Item	Approach	Avg Time	Time Std Dev
Screenshot	gRPC (<code>getScreenshot</code>)	19.94 ms	21.47 ms
View Hierarchy	ADB (<code>uiautomator</code>)	2.53 s	1.90 s
TOUCH action	gRPC (<code>sendTouch</code>)	419.50 μ s	64.71 μ s
LIFT action	gRPC (<code>sendTouch</code>)	412.30 μ s	84.18 μ s
TEXT action	ADB (<code>input</code>)	1.30 s	0.28 s

B.1 Mobile-Env simulator

The simulator of Mobile-Env defines the basic observation space and action space for GUI agent interaction. The observation includes both screenshot and VH (View Hierarchy) XML. The basic raw action space consists of three types of atomic actions, TOUCH, LIFT, and TEXT. To keep the interaction as efficient as possible, we adopt the gRPC¹⁶ interfaces of Android Emulator to obtain the screenshot and perform TOUCH and LIFT actions. As there are not any gRPC methods to obtain VH and perform TEXT, we have to adopt ADB (Android Debug Bridge) to implement them, which takes much longer time delay. A comparison of the efficiency of acquiring the observation components and performing actions is provided in Tab. 7.

B.2 Certificate unpinning

As described in § 3.2.1 in the main paper, certificate pinning prevents real-world apps from trusting the SSL certificates issued by the replay proxy. Thus, to establish truly controlled environments, we sought help from reverse engineering and developed three solutions. 1) Injecting the proxy’s certificates into a writable system image and making it mock a built-in system certificate. 2) Replacing the verifier used by the app at runtime via runtime instrument tool. 3) Repackaging the APK (installation package of the app) to remove its configurations regarding certificate pinning. We validated the solutions on a wide series of common-use apps and affirm that most apps can be dealt with by at least one solution.

B.3 Task state estimation

Mobile-Env leverages five primary types of signals to enable reliable task state estimation and produce intermediate instructions and rewards and episode ends accordingly at the right time.

Screen text Mobile-Env detects or recognizes some texts matching a specific pattern from the screenshot. An external OCR system is required.

Screen icon Mobile-Env detects or recognizes some icons matching a specific pattern from the screenshot. The pattern can be specified through an icon class or a reference image. An external icon recognition system is needed.

View hierarchy Mobile-Env checks if an expected node exists in VH of the current screen. VH is acquired through ADB `uiautomator`.

System log Mobile-Env monitors the system log for some lines matching a specific pattern. System log is acquired through ADB `logcat`.

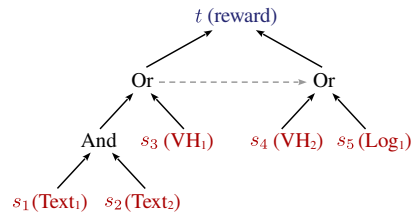
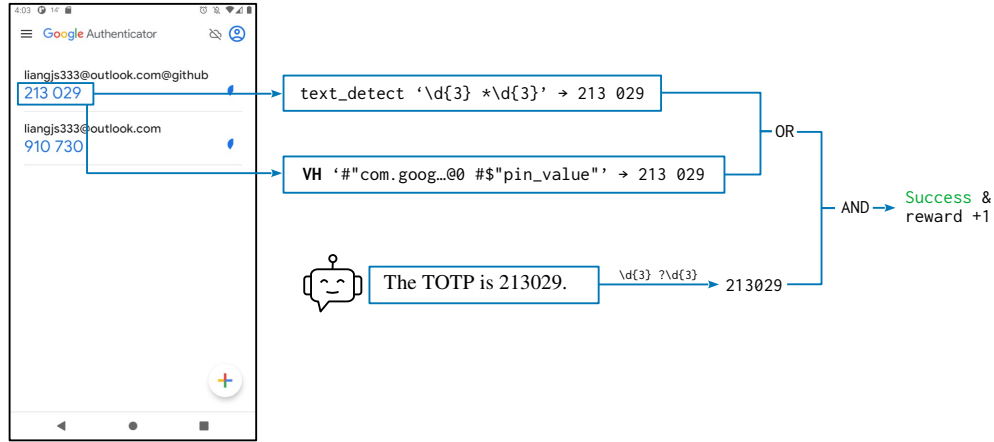
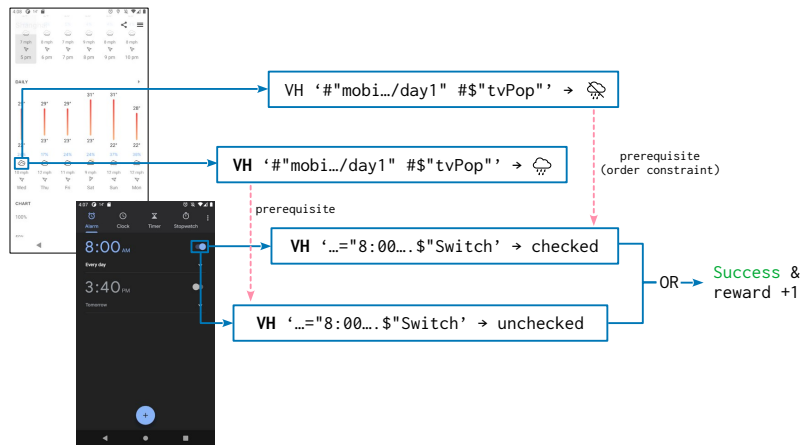


Figure 6: A demonstration of combining different types of system signals in a tree form. “Text”, “VH”, and “Log” indicate different types of signals.

¹⁶<https://grpc.io/>, an open-source general-purpose Remote Procedure Call (RPC) framework created by Google™.



(a) The instruction is “Check my TOTP for GitHub and tell me the code”. As the TOTP is a time-based dynamic passcode, the evaluator will compare the answer from the agent with the real-time passcode shown in the authenticator. The text is detected according to a regex (regular expression). The VH node is selected via a customized CSS selector (used to filter XML nodes).



(b) The instruction is “Turn off my alarm in the morning if it will rain tomorrow, otherwise leave it on”. The agent must first check the weather of tomorrow through the weather app on the phone before it switches the alarm clock on or off accordingly.

Figure 7: Examples of concrete task state estimation and evaluation via Mobile-Env

Response to human user (RHU) Mobile-Env listens to the response generated by the agent and checks if it matches the groundtruth answer.

Mobile-Env further introduces three strategies for combining these signals to enhance the flexibility of state estimation: conjunctive evaluation, disjunctive evaluation, and order constraint. Mobile-Env proposes to represent these combination strategies in a tree form. As shown in Fig. 6, we design AND and OR nodes to express conjunctive and disjunctive evaluations, respectively. Skip connections (the gray dashed line) are used to express the order constraints between different signals (events). In this form, the state estimation strategy can be expressed flexibly and clearly in the task configuration file. Two concrete examples of state estimation and reward generation via Mobile-Env are illustrated in Fig. 7.

B.4 Miscellaneous auxiliary tools

Annotation tool for human demonstration Touch and typing are close to human behaviors and are suitable to apply behavior cloning approaches. Based on Mobile-Env, an annotation tool is developed to collect human demonstrations for convenience. The annotation tool offers a straightforward web interface to collect the human actions. The annotators can directly click on the shown screenshot to mock touching on the mobile screen, or click on a button on the webpage to directly type a word from a shrunk vocabulary. The actions, observations, and task information will be recorded for behavior cloning.

Task definition template toolkit A simple suite of tools is developed to auto-generate task configuration files from templates. Such a task configuration template declares a group of value slots. The value slots can be instantiated according to a simple instantiation configuration file. Besides, the configuration instantiated from a single template can be combined sequentially to form a “multi-stage” task definition. The fixed world WikiHow task set is built with the help of this toolkit.

C Details about the task sets

C.1 Instructions of open world set

The full list of (the first) instruction in the open world set is listed in Tab. 8.

Table 8: The instructions of the open world set.

Instruction	Involved Apps	IR*	II†
I have multiple alarms to wake me up in the morning. Please check my alarms and tell me when I should get up finally tomorrow morning.	Alarm Clock	X	X
Help me to check-in in Qidian ¹⁷ app.	Qidian	X	X
Help me to open the “Do Not Disturb” mode for my phone.	System Settings	X	X
Turn off my alarm in the morning if it will rain tomorrow, otherwise leave it on.	Weather, Alarm Clock	X	X
Check my last email from GitHub	Gmail, Notein ¹⁸	X	✓
Please disable all the notifications of Gmail for me. I want a free weekend.	System UI	X	X
Send the illustration of Noelle I just downloaded to Lii Yihjiun through Gmail.	File Manager, Gmail	X	✓
Please help me to check the free spaces of my phone storage.	System Settings	X	X
Turn on my alarm in the afternoon.	Alarm Clock	X	X
Check my TOTP ¹⁹ for GitHub and tell me the code.	Google Authenticator	X	X
Turn off airplane mode for my phone.	System Settings	X	X
Turn on airplane mode for my phone.	System Settings	X	X
Tell me the Android version of my phone.	System Settings	X	X
Install open-source 2FA ²⁰ authenticator Aegis Authenticator from F-Droid ²¹ for my phone.	F-Droid	X	X
Install tool to connect my mobile with desktop KDE Connect from F-Droid for my phone.	F-Droid	X	X
Install MPD client on Android M.A.L.P. from F-Droid for my phone.	F-Droid	X	X
Install open-source bill management tool OpenMoneyBox from F-Droid for my phone.	F-Droid	X	X

* Intermediate Rewards † Intermediate Instructions

¹⁷A reading app

¹⁸A noting app

¹⁹Time-based One-Time Password

²⁰2-Factor Authentication

²¹An open-source app store

Table 8: The instructions of the open world set.

Instruction	Involved Apps	IR*	II†
Install famous “just works” video and music player VLC from F-Droid for my phone.	F-Droid	✗	✗
I have several funds and bills to be noted in OpenMoneyBox. Now open the app, enter the wizard, and help me to note the bills.	OpenMoneyBox	✓	✓
Uninstall the app Google Authenticator and install open-source 2FA authenticator Aegis Authenticator from F-Droid instead for my phone.	System UI, F-Droid	✓	✓
Uninstall the app VLC and install famous reimplementaion of classical MPlayer MPV from F-Droid instead for my phone.	System UI, F-Droid	✓	✓
Help me to set up a new alarm at 10:00 on Friday, Saturday, and Sunday.	Alarm Clock	✗	✗
Help me to set up a new alarm at 11:59 on Thursday and Friday.	Alarm Clock	✗	✗
Help me to set up a new alarm at 12:13 on Friday and Saturday.	Alarm Clock	✗	✗
Help me to set up a new alarm at 13:00 on Monday and Tuesday.	Alarm Clock	✗	✗
Help me to set up a new alarm at 14:34 on Thursday.	Alarm Clock	✗	✗
Help me to set up a new alarm at 15:00 on Wednesday, Thursday, and Saturday.	Alarm Clock	✗	✗
Help me to set up a new alarm at 16:00 on Tuesday.	Alarm Clock	✗	✗
Help me to set up a new alarm at 7:00 on Monday and Tuesday.	Alarm Clock	✗	✗
Help me to set up a new alarm at 8:00 on Friday.	Alarm Clock	✗	✗
Help me to set up a new alarm at 9:00 on Monday.	Alarm Clock	✗	✗
Uninstall the app KDE Connect on my phone.	System UI	✗	✗
Uninstall the app M.A.L.P. on my phone.	System UI	✗	✗
Uninstall the app OpenMoneyBox on my phone.	System UI	✗	✗
Search an article to learn how to bake lobster tails.	WikiHow	✓	✓
Search an article to learn how to deal with little sisters.	WikiHow	✓	✓
Search an article to learn how to determine shipping costs.	WikiHow	✓	✓
Search an article to learn how to file your nails.	WikiHow	✓	✓
Search an article to learn how to hide short bangs.	WikiHow	✓	✓
Search an article to learn how to how to behave yourself.	WikiHow	✓	✓
Search an article to learn how to how to find true happiness and peace.	WikiHow	✓	✓
Search an article to learn how to how to import from china into the usa.	WikiHow	✓	✓
Search an article to learn how to how to make a homemade weight set.	WikiHow	✓	✓
Search an article to learn how to make bottle cap art.	WikiHow	✓	✓
Search an article to learn how to make pupusas.	WikiHow	✓	✓
Search an article to learn how to rip dvd audio to mp3 using vlc media player.	WikiHow	✓	✓
Search an article to learn how to wear high waisted shorts.	WikiHow	✓	✓
Search an article to learn how to write a video game review.	WikiHow	✓	✓
Search an article to learn how to introduce science to preschoolers.	WikiHow	✓	✓
Search an article to learn how to make an electroscope.	WikiHow	✓	✓
Search an article to learn how to play the palace card game.	WikiHow	✓	✓
Search an article to learn how to rip jeggings.	WikiHow	✓	✓
Search an article to learn how to transfer a vehicle tag in florida.	WikiHow	✓	✓
Search an article to learn how to use scissors in gimp.	WikiHow	✓	✓

* Intermediate Rewards † Intermediate Instructions

Table 8: The instructions of the open world set.

Instruction	Involved Apps	IR*	II†
Search the interstellar.	IMDB	✓	✓
Navigate to the video part	IMDB	✓	✓
Search the first two topics about Hongloumeng in Douban	Douban	✓	✓
Get in the douban group of the book named Sanguoyanyi	Douban	✓	✓
Return the douban score of book named Sanguoyanyi	Douban	✓	✓
Tell me the 4th movie in the Douban Top250 movies list.	Douban	✓	✓
Search the Breaking Bad TV show on IMDB	IMDB	✓	✓
Search the 2002 verison spiderman movie and check the rate for it.	IMDB	✓	✓
Search the book named Sanguoyanyi.	Douban	✓	✓
Return the director of the TV named Zhenhuanzhuan	Douban	✓	✓
Open the address table.	Ele.me	✓	✓
Help me order takeout.	Ele.me	✓	✓
Help me order takeout.	Ele.me	✓	✓
Select one of my previous orders from my entire order history and place a new order of the same items to help me order takeout again.	Ele.me	✓	✓
Help me order a cup of my favorite milk tea by searching for it.	Ele.me	✓	✓
Help me order a cup of my favorite milk tea by searching for Coco (the store name).	Ele.me	✓	✓
Help me select takeout that can be delivered within 30 minutes using the filters on the homepage and place the order.	Ele.me	✓	✓
View my order history.	Ele.me	✓	✓
View my shopping cart.	Ele.me	✓	✓
View the messages page.	Ele.me	✓	✓

* Intermediate Rewards † Intermediate Instructions

C.2 App data crawling and replay for WikiHow task set

In order to establish a truly reproducible and comparable benchmark, we need to crawl and store the varying app data of WikiHow to build the fixed world set. The app data will be replayed at runtime to mimic the real scenario.

WikiHow app comprises two types of pages: static pages and dynamic pages. The static pages are referred to as the pages independent of user inputs²² such as the home page, articles, author information, *etc.* The dynamic pages are the pages sensitive to user inputs, such as the pages of search results. App data for static pages are directly collected from WikiHow website²³. As observed that the pages in the app are different from the pages in the browser, we checked the HTTP request headers of all the types of pages and mocked the app requests to obtain the correct contents. Considering that the user always explores the app from the home page, the crawler also starts from the home page and accesses the links in breadth-first order. For simplicity, the links to the outside of WikiHow are omitted. All the links including all types of media are requested and stored for a perfect replay. Finally, 856,045 resources are dumped, in which there are 107,448 distinct pages. These website resources totally occupy about 88 GiB. In contrast to the static resources, it is nearly impossible to traverse all the probable combinations of keywords to dump the pages of search results. Consequently, an open-source information retrieval engine, Pyserini [9], is adopted following Yao et al. [31] to mimic the WikiHow search engine. Then the synthetic search pages are built dynamically during the replay upon the template of the search result page in WikiHow app.

²²Static pages may contain contents that vary according to spatio-temporal conditions as well, *e.g.*, dynamic recommendations, statistics for visits, and ratings. However, these contents do not depend on instant user inputs, and thus are considered “static” resources.

²³<https://www.wikihow.com/Main-Page>

The crawled data are tagged with the requested URL paths. Each response is dumped in an individual file. Both HTTP headers and raw payload are saved. The flow data can be replayed through mitmproxy²⁴. The saved static resources are directly read and replayed with the timestamps refreshed. The search result page is crafted with the results from Pyserini engine [9] and the page template. Pyserini builds the indices of the articles from the article titles. The page template is built from a real search result page in WikiHow app by making the search result list a filling slot.

C.3 Automatic task generation for WikiHow task set

By thoroughly exploring WikiHow app, 10 basic operations on WikiHow are summarized as:

- Search for an article,
- Access an article from the current page,
- Access the about page from article page,
- Access author information page from article page,
- Access category content page from article page,
- Check reference list,
- Rate yes/no for article,
- Share article,
- Bookmark article,
- Check the bookmarks.

Based on this summary, a suite of “single-stage” sub-task configuration templates is designed:

- `home2search`, instructing to search for an article from the home page.
- `search2article`, `author2article`, & `category2article`, instructing to access an article from search result page, author information page, and category content page, respectively.
- `article2about`, instructing to access the about page from article page.
- `article2author`, instructing to access author information page from article page.
- `article2category`, instructing to access category content page from article page.
- `article2reference`, instructing to check reference list on article page.
- `article2rate_no`, instructing to rate no for article
- `article2rate_yes`, instructing to rate yes for article
- `article2share`, instructing to share article
- `article2bookmark`, instructing to bookmark article and then check the bookmarks.

among which, `article2reference`, `article2rate_no`, `article2rate_yes`, `article2share`, and `article2bookmark` are used to form in-page tasks. In addition, we explore two open-source QA datasets based on WikiHow, `wikihow-lists`²⁵ and `WikiHowNFQA`²⁶, to construct QA tasks in WikiHow task set. `wikihow-lists` comprises three types of question:

- `stepped_summary`, expecting a list of key steps summarized from the current article.
- `ingredients`, expecting a list of ingredients for a recipe extracted from the current article.
- `needed_items`, expecting a list of needed items for an activity extracted from the current article.

`WikiHowNFQA` posts free-form summarization questions related to WikiHow articles. Based on these question-answer pairs, we additionally craft 4 sub-task configuration templates:

- `article2steps`, crafted from `stepped_summarys`

²⁴<https://mitmproxy.org/>

²⁵https://huggingface.co/datasets/b-mc2/wikihow_lists

²⁶<https://huggingface.co/datasets/Lurunchik/WikiHowNFQA>

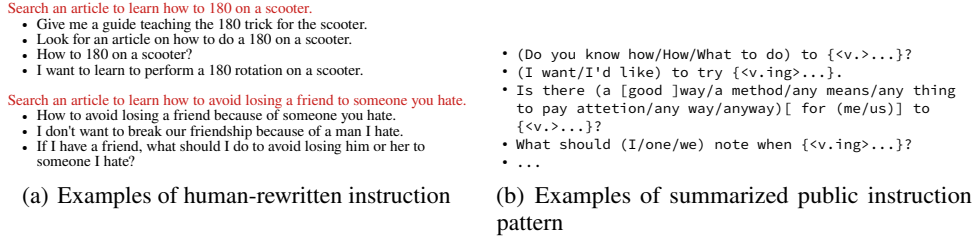


Figure 8: Examples of human-rewritten instructions (Left) and public instruction patterns (Right)

- `article2ingredients`, crafted from `ingredientss`
- `article2needed_items`, crafted from `needed_items`
- `article2summary`, crafted from WikiHowNFQA questions

Each template contains a group of filling slots expecting some keywords like article title, author name, question, and groundtruth answer. Then these keywords are sampled from the crawled app data or from the two QA datasets to instantiate the templates. Subsequently, the instantiated templates are concatenated into multi-stage task definitions under the constraint that the target page/element/answer (the part after “2”, *e.g.*, “share” from “article2share”) is directly on/referenced by the current page (the part before “2”, *e.g.*, “article” from “article2share”). Finally, we obtained the task set of 150 multistage tasks in which there are 2.68 single-stage sub-tasks averagely.

C.4 Instruction rewriting & dynamic replacement for WikiHow task set

The template-generated task definitions of WikiHow are homogeneous and straightforward, directly containing the target article name to search or access. These instructions are too simple for LLMs to understand. In order to establish a more challenging benchmark and evaluate the capacity of LLM for instruction understanding, the authors invited several colleagues and leveraged ChatGPT to rewrite the stereotyped instructions. Each person created 2–3 new expressions for each instruction. To further improve diversity, the authors inspected all the rewritten instructions and summarized several public expression patterns that are easy to transfer to different topics. Several examples of rewritten instruction and public pattern are depicted in Fig. 8. The corresponding rewritten instructions and public patterns are randomly sampled to replace the original ones at runtime. The different branches in pattern (*e.g.*, (I want/I'd like)) are randomly selected and the keyword in slots (*e.g.*, {<v.ing>...}) are transformed into the correct POS (part-of-speech). POS recognition and transformation are implemented with the help of NLTK²⁷ and lemminflect²⁸.

D Details about LLM-based agents and experiments

D.1 LLM access and compute resources

We tested the closed-source text LLMs including GPT-3.5 (gpt-3.5-turbo-1106 & gpt-3.5-turbo-instruct), GPT-4²⁹ (gpt-4-1106-preview), and Claude-3³⁰ (claude-3-opus-20240229). The capable closed-source VLMs are also tested: GPT-4V (gpt-4-1106-vision-preview) and Claude-3 (claude-3-opus-20240229). These models are accessed through the online APIs (Application Programming Interfaces) provided by OpenAI and Anthropic.

We tested open-source LLMs including LLaMA 2 [22] (llama-2-chat-7b, llama-2-chat-13b, and llama-2-chat-70b), AgentLM [33] (agentlm-7b, agentlm-13b, and agentlm-70b), and LLaMA 3 [1] (llama-3-8b-instruct & llama-3-70b-instruct). These models are deployed on NVIDIA A10 Tensor

²⁷<https://www.nltk.org>

²⁸<https://github.com/bjasacob/LemmInflect>

²⁹<https://platform.openai.com>

³⁰<https://www.anthropic.com/claude>

Table 9: Mapping from VH node classes to HTML tags

VH node class	HTML tag class
*TextView	p
*Button	button
*MenuItemView	button
*ImageView	img
*IconView	img
*Image	img
*EditText	input (type="text")
others	div

Table 10: Mapping from VH properties to HTML attributes. The resource-id property usually consists of three parts: the package name, the element class, and the element name. Only the element name is used for the converted class attribute, which follows Wang et al. [24].

VH property	HTML attribute
resource-id	class
content-desc	alt

```
<hierarchy rotation="0">
  <node index="0" text="" resource-id="" class="" ...>
    <node index="0" text="" resource-id="" class="" ...>
      <node index="0" text="" resource-id="" class="" ...>
        <node index="0" text="" resource-id="" class="" ...>
          ...
        <node index="0" text="" resource-id="" class="" ...>
          ...
        <node index="0" text="" resource-id="" class="" ...>
          <node index="0" text="" resource-id="" class="" ...>
            <node index="0" text="" resource-id="" class="" ...>
              <node index="0" text="" resource-id="" class="" ...>
                <node index="1" text="" resource-id="" class="" ...>
                  <node index="2" text="" resource-id="" class="" ...>
                    ...
                  </node>
                </node>
              <node index="1" text="" resource-id="" class="" ...>
                <node index="0" text="" resource-id="" class="" ...>
                  ...
                </node>
              </node>
            </node>
          </node>
        </node>
      </node>
    </node>
  </node>
</hierarchy>
```

(a) Raw view hierarchy

```
<node class="android.widget.ImageButton" resource-id=""
content-desc="Open navigation drawer" text=""
clickable="true" id="0"/>
<node class="android.widget.ImageView" resource-id=
"com.wikihow.wikihowapp:id/wikihow_toolbar_logo"
content-desc="" text="" clickable="false" id="1"/>
<node class="android.widget.ImageView" resource-id=
"com.wikihow.wikihowapp:id/search_button"
content-desc="Search" text="" clickable="true"
id="2"/>
<node class="android.webkit.WebView" resource-id=
"com.wikihow.wikihowapp:id/webView" content-desc=""
text="" clickable="true" id="3"/>
<node class="android.view.View" resource-id=
"android:id/statusBarBackground" content-desc=""
text="" clickable="false" id="4"/>
```

(b) Visible VH leaves (simplified)

```
<button alt="Open navigation drawer" id="0" clickable=
"true"></button>
<img class="wikihow toolbar logo" id="1" clickable=
"false">
<img class="search button" alt="Search" id="2" clickable=
"true">
<div class="webView" id="3" clickable="true"></div>
<div class="statusBarBackground" id="4" clickable=
"false"></div>
```

(c) Converted HTML elements

Figure 9: Example of conversion from VH into HTML. Fig. (a) depicts the raw view hierarchy extracted from the OS. Fig. (b) depicts the visible leaf nodes from Fig. (a). Additional id property is appended. Fig. (c) depicts the converted HTML elements from Fig. (b) according to the conversion rule.

Core clusters with bfloat16 with the help of Hugging Face Text-Generation-Inference³¹ and vLLM³². 7B (8B), 13B, and 70B models require 1, 2, and 8 A10 GPUs with 24 GiB memory, respectively.

In our experiments, Android Emulator runs on a Linux host with KVM (Kernel-based Virtual Machine) enabled to conduct efficient enough emulation.

D.2 Observations for LLM/VLM-based agents

The text LLMs are fed with a simplified HTML format derived from the VH component in the raw observation following Wang et al. [24]. In the HTML, only the visible leaf nodes in VH are preserved. The VH classes and essential properties are converted into HTML tags and attributes. Tab. 9 depicts the mapping rule from classes of VH nodes to HTML tags. Tab. 10 depicts how the

³¹<https://github.com/huggingface/text-generation-inference>

³²<https://docs.vllm.ai/en/latest/index.html>

Table 11: Actions spaces for different types of LLMs/VLMs

Model Type	Action	Description
Text LLM	CLICK(element_id)	Click on an element.
	LONG_CLICK(element_id)	Long click on an element.
	INPUT(element_id, text)	Click on an element and input texts into it.
	SCROLL(direction)	Scroll the page (<i>i.e.</i> slide towards the opposite direction).
	ANSWER(text)	Generate an answer response to human user.
	GOBACK	Go back to the last page.
VLM w/ SoM	TAP(element_id)	Tap an element.
	LONG_TAP(element_id)	Long tap an element.
	TYPE(text)	Directly type some texts.
	SLIDE(direction)	Slide the “finger” towards a direction.
	ANSWER(text)	Generate an answer response to human user.
	GOBACK	Go back to the last page.
VLM w/o SoM	TAP(x, y)	Tap a pixel.
	LONG_TAP(x, y)	Long tap a pixel.
	TYPE(text)	Directly type some texts.
	SLIDE(x0, y0, x1, y1)	Slide the “finger” from (x0, y0) to (x1, y1).
	ANSWER(text)	Generate an answer response to human user.
	GOBACK	Go back to the last page.

VH node properties are mapped to the HTML element attributes. An example of conversion from VH into HTML is illustrated in Fig. 9.

The VLMs are fed with screenshots applied Set-of-Marks (SoM) [29]. We adopt EasyOCR³³ to detect the screen texts and vision-ui³⁴ to detect the screen icons and mark the detected elements using numeric tags.

D.3 Action spaces for LLM/VLM-based agents

The designed action spaces for models are shown in Tab. 11.

These actions (*i.e.*, CLICK, LONG_CLICK, INPUT, SROLL, TAP, LONG_TAP, TYPE, and SLIDE) are supposed to be translated into basic actions (*i.e.*, TOUCH, LIFT, and TEXT) to be executed in Mobile-Env. Particularly, CLICK, LONG_CLICK, TAP, and LONG_TAP are translated into a sequence of TOUCHs appended with a LIFT, *e.g.*,

$$\text{CLICK}(eid) \mapsto \underbrace{\text{TOUCH}(c), \dots, \text{TOUCH}(c)}_{n \text{ TOUCHs}}, \text{LIFT}. \quad (1)$$

Here eid is the integer element id and c is the center point of the element bounding box. n is set to 3 for CLICK and TAP and 10 for LONG_CLICK and LONG_TAP in our experiments. TYPE is translated into a sequence of TEXTs, and a keyboard action “Enter”:

$$\text{TYPE}(\text{text}) \mapsto \text{TEXT}(\text{tkn}_1), \dots, \text{TEXT}(\text{tkn}_l), \text{Enter}, \quad (2)$$

where tkn_i denotes the i -th token from the input text . The keyboard action “Enter” is a special action beyond the basic actions and indicates typing a “carriage return” into the OS. INPUT is translated into a CLICK followed by a TYPE. SCROLL and SLIDE are translated into m consecutive TOUCHs followed by a LIFT, where m depends on the sliding distance. These TOUCHs form a consecutive sliding of the “finger” on the screen. Fig. 10 demonstrates the translated sliding trajectories of typical SCROLLs.

D.4 Hyperparameter selection

It is found that the maximum length of 2 exemplars can just fit in the 4K-token context of the open-source models (LLaMA 2, AgentLM, and LLaMA 3). Thus, 2-shot in-context learning (ICL)

³³<https://www.jaided.ai/easyocr/>

³⁴<https://github.com/Meituan-Dianping/vision-ui>

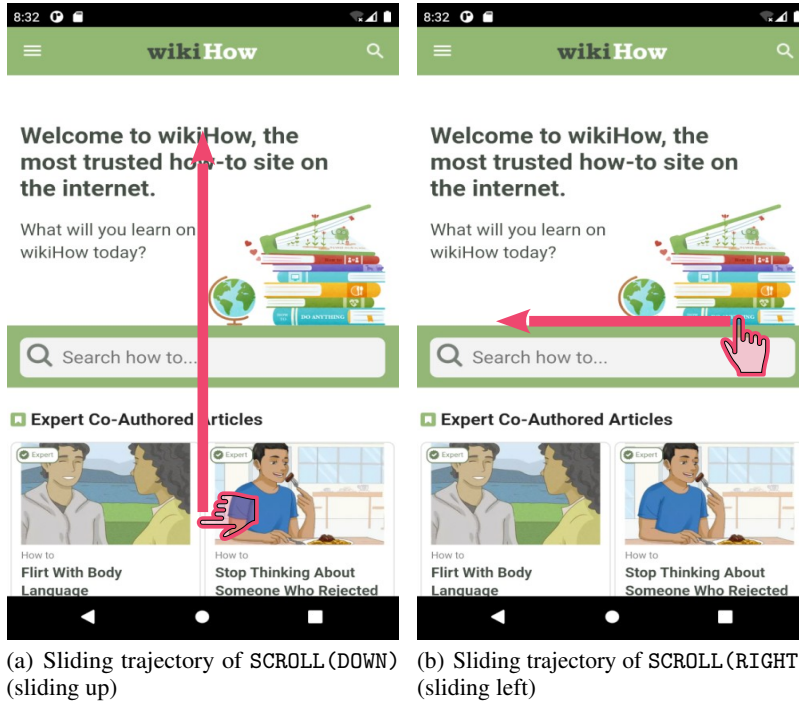


Figure 10: Demonstration of sliding trajectories translated from SCROLL. The starting point is at 20% of the screen height/width, while the terminating point is at 80%.

is adopted for a fair experiment setting. A relative low sampling temperature 0.1 is selected to balance the generation of structured action representations and free-form reasoning languages, also to balance the reproducibility and exploration ability of models. Before formal experiments, the authors invited several colleagues to complete the designed tasks. This is implemented by simply replacing the revocation of LLM with keyboard inputting. Text observations are simply printed to the terminal, and screenshot observations are shown in a separate window. Their feedback reveals that the aforementioned observation & action spaces are sufficient to complete the tasks, and most of the tasks can end in 5 steps. Therefore, the episodes with more than 15 steps are considered failures during formal evaluation.

D.5 Prompts for LLM/VLM-based agents

The complete prompts for LLMs and VLMs are provided in Tab. 12 and Tab. 13.

Table 12: Prompt for text LLMs

System:	Given several information about the task and the screen at the current step, I need to take an appropriate action according to the given information to finish the task. Available actions are:
	<pre> INPUT(element_id, text) # You can input something into search area through this action CLICK(element_id) # You can click on some clickable element through this action LONG_CLICK(element_id) # You can long lick on some clickable element through this action SCROLL(direction) # You can scroll up/down to browse the long page through this action </pre>

Table 12: Prompt for text LLMs

	ANSWER(text) # You can generate an answer to me through this action
	GOBACK # You can return to the last page through this action
	At every time decision, you need first output a line of think, and then output a line of action. For example:
	THINK: I need to have a thinking before I take my action.
	ACTION: ANSWER(I can take any available action, e.g., give an answer.)
	The above is just an example and takes no any relation with the following inputs. Now, take your mission.
User:	<pre> ''' Task: Look up how to conceal gauges. Then, access the article regarding gauges. Screen: <button alt="Open navigation drawer" id="0" clickable="true"></button> <div class="webView" id="3" clickable="true"></div> <div class="statusBarBackground" id="4" clickable="false"></div> Instruction: Action History: ---</pre>
Assistant:	<pre> THINK: I should search "hide gauges" to know how to conceal gauges. ACTION: INPUT(2, hide gauges) '''</pre>
User:	<pre> Task: How to style ruby rose hair? Then, open the article regarding ruby. Then, check the about page for the reliability of WikiHow. Screen: <button alt="Open navigation drawer" id="0" clickable="true"></button> <input class="search src text" value="Do ruby rose hair " type="text" id="1" clickable="true"> <div id="3" clickable="true">How to Do Ruby Rose Hair</div> <div id="4" clickable="true">● </div> <p id="5" clickable="true">41,446 views</p> <div id="6" clickable="true">● </div> <p id="7" clickable="true">Updated</p> <p id="8" clickable="true">2 years ago</p> <div id="9" clickable="true">● </div> <p id="10" clickable="true">Expert Co-Authored</p> <div id="11" clickable="true">How to Dye Your Hair Rose Gold</div> <div id="12" clickable="true">● </div> <p id="13" clickable="true">48,548 views</p></pre>

Table 12: Prompt for text LLMs

```

<div id="14" clickable="true">● </div>
<p id="15" clickable="true">Updated</p>
<p id="16" clickable="true">3 years ago</p>
<div id="17" clickable="true">● </div>
<p id="18" clickable="true">Expert Co-Authored</p>
<div id="19" clickable="true">How to Do a Hair Mask for Frizzy
Hair</div>
<div id="20" clickable="true">● </div>
<p id="21" clickable="true">29,415 views</p>
<div id="22" clickable="true">● </div>
<p id="23" clickable="true">Updated</p>
<p id="24" clickable="true">2 years ago</p>
<div id="25" clickable="true">● </div>
<p id="26" clickable="true">Expert Co-Authored</p>
<div id="27" clickable="true">How to Do 1920s Hair</div>
<div id="28" clickable="true">● </div>
<p id="29" clickable="true">194,834 views</p>
<div id="30" clickable="true">● </div>
<p id="31" clickable="true">Updated</p>
<p id="32" clickable="true">1 year ago</p>
<div id="33" clickable="true">● </div>
<p id="34" clickable="true">Expert Co-Authored</p>
<div id="35" clickable="true">How to Do 50s Hair</div>
<div id="36" clickable="true">● </div>
<p id="37" clickable="true">39,890 views</p>
<div id="38" clickable="true">● </div>
<p id="39" clickable="true">Updated</p>
<p id="40" clickable="true">1 year ago</p>
<div id="41" clickable="true">● </div>
<p id="42" clickable="true">Expert Co-Authored</p>
<div id="43" clickable="true">How to Do 70s Hair</div>
<div id="44" clickable="true">● </div>
<p id="45" clickable="true">69,458 views</p>
<div id="46" clickable="true">● </div>
<p id="47" clickable="true">Updated</p>
<p id="48" clickable="true">1 year ago</p>
<div id="49" clickable="true">● </div>
<p id="50" clickable="true">Expert Co-Authored</p>
<div id="51" clickable="true">How to Do Goku Hair</div>
<div id="52" clickable="true">● </div>
<p id="53" clickable="true">44,049 views</p>
<div class="statusBarBackground" id="54"
clickable="false"></div>

```

Instruction:

"How to Do Ruby Rose Hair" may be helpful for us.

Action History:

INPUT(2, do ruby rose hair) # I should search "do ruby rose hair" to know how to style ruby rose hair.

Assistant: THINK: The instruction requires to read the article "How to Do Ruby Rose Hair". Element 3 is "How to Do Ruby Rose Hair" and is clickable. I should click it to access the article.

ACTION: CLICK(3)

Table 13: Prompt for VLMs

System:	<p>Given several information about the task and the screen at the current step, I need to take an appropriate action according to the given information to finish the task. Available actions are:</p> <p>TAP(element_id) # You can touch a marked element on the screen by its number id. LONG_TAP(element_id) # You can long touch a marked element on the screen by its number id. SCROLL(direction) # You can scroll up/down to browse the long page through this action TYPE(text) # You can type in something through keyboard. ANSWER(text) # You can generate an answer to me through this action. GOBACK # You can return to the last page through this action.</p> <p>At every time decision, you need first output a line of think, and then output a line of action. For example:</p> <p>THINK: I need to have a thinking before I take my action. ACTION: ANSWER(I can take any available action, e.g., give an answer.)</p> <p>The above is just an example and takes no any relation with the following inputs. Now, take your mission.</p>
User:	<p>Task: Look up how to conceal gauges. Then, access the article regarding gauges. Instruction:</p>
Assistant:	<p>Action History: TAP(13) # I should search "hide gauges" to know how to conceal gauges. I can tap on the search box to input the search keywords. Screen: <code>#{imagef:prompts/mm_prompt/eg1_som.png:280:500}</code> THINK: I should search "hide gauges" to know how to conceal gauges. I have selected the search box and need to type the search keywords. ACTION: TYPE(hige gauges)</p>
User:	<p>Task: How to style ruby rose hair? Then, open the article regarding ruby. Then, check the about page for the reliability of WikiHow. Instruction: "How to Do Ruby Rose Hair" may be helpful for us. Action History: TAP(13) # I should search "Do ruby rose hair" to know how to style ruby rose hair. I can tap on the search box to input the search keywords. TYPE(Do ruby rose hair) # I should search "Do ruby rose hair" to know how to style ruby rose hair. I have selected the search box and need to type the search keywords. Screen: <code>#{imagef:prompts/mm_prompt/eg2_som.png:280:500}</code></p>
Assistant:	<p>THINK: The instruction requires to read the article "How to Do Ruby Rose Hair". Thus, I should click on this link. ACTION: TAP(3)</p>

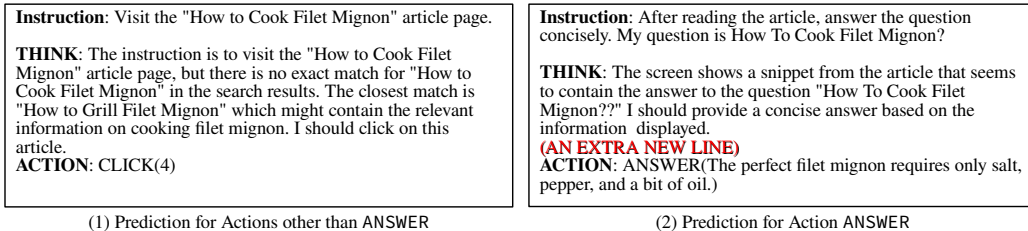


Figure 11: Different behavior of GPT-4 when predicting different types of action. When outputting ANSWER action, the model is prone to predict an extra “new line” between the “THOUGHT” and “ACTION” lines.

D.6 More analyses on open world set

Agents fail to open correct apps from the Android homepage.

To reduce difficulty, the environments are launched with the involved app open before the agent’s execution. Wondering the raw performance of the agents, we further choose the first 34 tasks in Tab. 8 (the tasks involving various apps beyond WikiHow, Douban, IMDB, and Ele.me) and conduct ablation study by making the agent start the task from the homepage. The results are shown in Tab. 14. In this setting, the involved app is not open in advance, and the agent has to open the correct app by itself. This seems to be hard for most models. The difficulty mainly lies on the way to open the app menu. The correct approach is to slide up on the home page, but the agents often fail to predict this action. The agents may also be misled, *e.g.*, they always mistake the icon of Google (``) for the icon of Gmail and approach failure.

Table 14: Ablation study regarding the difficulty of “starting from homepage”. “App” indicates the involved app is open at the task beginning, while “Home” indicates the app is not open and the agent has to complete the task by open the correct app firstly. However, the agents struggle to open the correct app when facing the Android home page, and thus obtain worse results.

Model	SR		Rwd	
	App	Home	App	Home
Claude-3-Opus	14.70	8.83	0.18	0.09
GPT-4	38.23	23.53	0.41	0.26
AgentLM-70B	8.82	0.00	0.09	0.00
LLaMA-3-70B	20.59	12.50	0.21	0.15

There are other difficulties revealed by our experiments that hinder the LLM and VLM agents from completing tasks. For example, VH is not available on several sensitive pages, *e.g.*, the page showing the version of Android OS. Thus, such tasks can only be completed by visual models like GPT-4V. The safety barrier of the models can also lead to failed execution, *e.g.*, Claude-3-Opus refuses to check the TOTP for the human user. Inexplicable strange activities also exist, such as GPT-4V always existing the WikiHow app at once when encountering tasks on WikiHow.


D.7 More analyses & case study on fixed world set

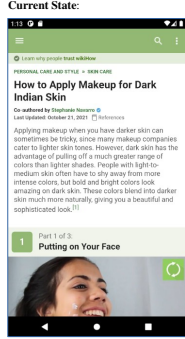
Failure of GPT-4 on QA tasks GPT-4 mysteriously predicts an extra “new line” between its thought and ANSWER action, which violates the output format and leads to failure. A group of examples is shown in Fig. 11. As discussed in § 5.2 in the main paper, this phenomenon has become a significant bottleneck of the capable models’ performance.

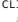
Capability difference of chat and instruct models The chat model (GPT-3.5) and instruct model (GPT-3.5-instruct) do not show significant disparity in the average performance (37.33% SR vs. 36.67% SR), however, demonstrate the opposite advantages of completing in-page and cross-page tasks (38.98% SR on cross-page tasks and 49.02% SR on in-page tasks of GPT-3.5 & 44.07% SR on cross-page tasks and 35.29% SR on in-page tasks of GPT-3.5-instruct). Looking through the relative task samples, two typical cases are depicted in Fig. 12. It seems that the instruct model is more prone to ignore the most recent instruction and try to achieve a stale goal. The chat model is more adept at finding a target link through scrolling and attempting (maybe a wrong link). However, sometimes,

Task Description:
I need to know some method to apply makeup for dark indian skin. Then, visit the article page with regard to "skin". Then, make a bookmark for the article and then check it.

Current Instruction:
After bookmarking the article, go check the bookmark in the bookmark list.

Action History:
INPUT (G, how to apply makeup for dark Indian skin)
CLICK("How to Apply Makeup for Dark Indian Skin")
CLICK()
CLICK("Bookmark")



Chat Model:
CLICK() # The instruction is to add a bookmark for the article and then go to the bookmark list to check. →
CLICK("Bookmarks") # The instruction requires me to add a bookmark for the current article "How to Apply Makeup for Dark Indian Skin" and then go to the bookmark list to check. →
SUCCESS

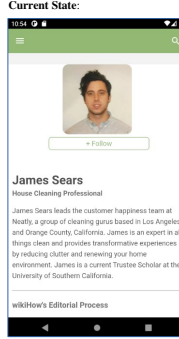
Instruct Model:
INPUT("SKIN CARE", apply makeup for dark indian skin) # I should search "apply makeup for dark indian skin" to find an article about it. Element 7 is the search bar and is clickable. I should input "apply makeup for dark indian skin" into it.

(a) Chat model outperforms instruct model on in-page tasks.

Task Description:
How should I brush my nonstick pan after cooking? Then, open the article about "clean", "pan". Then, check out James's author page. Then, the target page is involving "vinegar".

Current Instruction:
Open the article "How to Clean a Keurig with Vinegar".

Action History:
INPUT (G, clean nonstick pan)
CLICK("How to Clean a Nonstick Pan")
CLICK("James Sears")



Chat Model:
CLICK(11) # →
SCROLL(DOWN) # The instruction is to visit the article page "How to Clean a Keurig with Vinegar". I need to find the link to this article. →
SCROLL(DOWN) # The instruction is to visit the article page "How to Clean a Keurig with Vinegar". I need to scroll down to find the article link. →
...

Instruct Model:
INPUT(G, how to clean a keurig with vinegar) # I should search "how to clean a keurig with vinegar" to access the article. →
CLICK("How to Clean a Keurig with Vinegar") # The instruction requires to read the article "How to Clean a Keurig with Vinegar". Element 3 is "How to Clean a Keurig with Vinegar" and is clickable. I should click it to access the article. →
SUCCESS

(b) Instruct model performs better on cross-page tasks.

Figure 12: Capability difference of chat model and instruct model. For clarity, the element ids are replaced with its text or icon in the figure. Non-existing element id is marked with red. Thought texts are rendered in gray and separated with action text by "#".

Table 15: Results of AgentLMs on WikiHow task set

Model	Rwd	SR	Cross-Page		In-Page		QA	
			Rwd	SR	Rwd	SR	Rwd	SR
AgentLM-7B	1.39	12.67	1.24	23.73	1.61	9.80	1.34	0.00
AgentLM-13B	1.44	14.67	1.31	28.31	1.65	5.88	1.38	5.00
AgentLM-70B	1.71	30.00	1.75	44.07	1.80	19.61	1.55	22.50

it requires remarkable scrolling down steps to achieve the goal. In contrast, the instruct model can make use of search again to quickly finish the new instruction.

Performances of AgentLMs The results of AgentLMs on the closed world set are depicted in Tab. 15. AgentLM models are further fine-tuned based on LLaMA-2 chat models through AgentTuning to adapt to decision-making tasks. They are expected to perform much better than the original common-use models. And the results do prove their effectiveness. AgentLM models gain a remarkable improvement on WikiHow task set compared to the original LLaMA-2 models, especially the 7B and 70B models. The strength of AgentLMs mainly lie in the more capable ability to follow the output format, *i.e.*, less often outputting illegal actions like SHARE and ANSWER: . . . (the correct format is ANSWER(. . .)). This fact is proven by the statistics in Tab. 16.

Failure of AgentLM-7B on QA tasks AgentLM-7B fails completely on QA tasks. The trajectories (cases shown in Fig. 13) reveal that AgentLM-7B often violates the output protocol and predicts invalid actions. Sometimes, it also repeats the certain question as the answer, even a proper summary has already been included in the thought. These phenomena reveal that AgentLM-7B is still weak at format following.

Results on single-stage sub-tasks Fig. 14 visualizes the performance of the models on the single-stage sub-tasks (introduced in § C.3). It is obvious that almost all the models fail on sub-tasks article2rate_no, article2rate_yes, author2article, and category2article. This is caused by the need of long-term scrolling down to complete these tasks. The buttons to give a rate for the article is at the page bottom. To find the target article from the author information page or

Table 16: Ratios of invalid actions predicted by models. AgentLMs predict much less invalid actions than models from LLaMA 2.

	7B	13B	70B
LLaMA-2	81.27	33.99	29.87
AgentLM	15.44	13.45	3.12

THINK: The article "How to Check if a Remote Control is Transmitting an Infrared Signal" seems to be relevant to the task. I should read through the article to find the needed information.
ACTION: READ(4)

THINK: The article "How to Check if a Remote Control is Transmitting an Infrared Signal" provides the needed information. I should encapsulate the important points.
ACTION: SUMMARIZE(4)

THINK: THINK: The article "How to Cook Filet Mignon" provides a step-by-step guide on how to cook filet mignon. The first step is to season the filet mignon with salt, pepper, and garlic. Then, ...
ACTION: ANSWER(How To Cook Filet Mignon?)

Figure 13: AgentLM-7B hallucinates invalid action types or returns the question itself as the final answer.

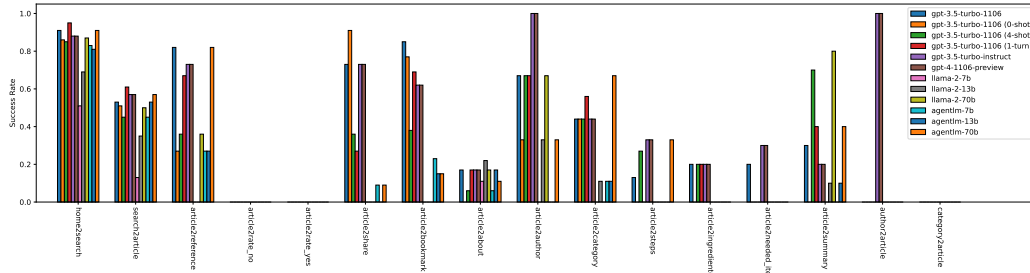
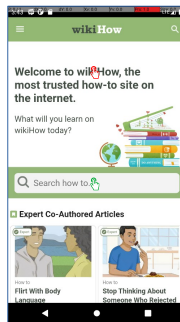
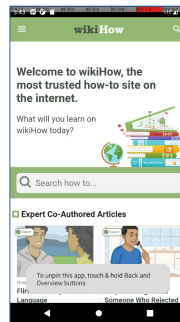


Figure 14: Success rates on single-stage sub-tasks

category content page also requires scrolling down for a number of steps. However, the agents may lack the belief to continuously scrolling down. They will try to click something aiming at reaching the target quickly during scrolling, which results in an inefficient policy and a failure after reaching the step number limit (*i.e.*, 15). An example of GPT-3.5 is given in Fig. 15. There are also shortcuts to complete these tasks more efficiently, *e.g.*, for `article2rate_no` and `article2rate_yes`, the link to the reference list will bring the agent to the very bottom of the article page near the rating buttons, and for `author2article` and `category2article`, the agent can find the target article through searching rather than seeking for it from the current page. However, this may require some *a priori* common-sense knowledge, and few models can leverage this. If a learnable agent conducts some prior explorations before testing (*e.g.*, reinforcement learning training, or building a grounded world model), it may leverage this *a priori* knowledge much better. This point can be studied in future work.



Task Description: I want to find a good way to bypass youtube's regional filter. Then, open the article about "YouTube". Then, enter the category Computers and Electronics. Then, the target article is regarding "play", "Amazon".
Current Instruction: Do you know how to bypass youtube's regional filter?
Agent: TAP(0.5, 0.2) # I need to search for information on how to get around YouTube's regional filters. I should tap on the search box to input the search keywords.
Proper Position: TAP(0.5, 0.55)



Agent: ANSWER(To unpin this app, touch & hold Back and Overview buttons) # The screen is still not displaying the content properly after multiple attempts to go back. The usual method of going back is not resolving the issue, and I do not have the capability to restart the app or check the device's internet connection. However, since the screen is showing a message about unpinned apps, it seems like the app might be pinned. I should try to unpin the app as suggested by the on-screen message.
GOBACK: # The screen is still not displaying the content properly after multiple attempts to go back. The previous action suggested that the app might be pinned. Following the on-screen message's advice to unpin the app could be the next logical step, but since I can't perform that action, I should try to go back again as it's the only available action I can take.

(a) VLM-based agent fails to ground the action accurately. (b) VLM-based agent is distracted by the popup.

Figure 16: Failure cases of VLMs

Failure cases of VLMs The VLM-based agent without SoM has a poor capability to ground the action. As shown in Fig. 16(a), even though the VLM predicts the plan correctly, it failed to accurately predicts the coordinates. Meanwhile, VLM is more prone to be distracted by some notification popups that are not present in the VH. For instance, in Fig. 16(b), the VLM is distracted by the popup about screen pinning and turns to try to unpin the screen rather than focuses on completing the task goal.

E Limitations

Despite the contribution of this work, there are still several limitations. One of the concern is about data contamination of WikiHow task set. We should claim that even though WikiHow is supposed to be a preferred corpus for LLM pretraining, the involved data formats in WikiHow task set, regardless of text-based or visual-based observations, are quite different from those that will be used in pretraining. And the problem formulation also deviates from the pretraining tasks of most common-use LLMs. However, there will be a latent risk for QA tasks, *i.e.*, the LLM may answer the question relying on its parametric memory, rather than according to the underlying article. Currently, we haven't worked out an efficient solution to detect such contamination, and we plan to leave it to future work for more efficient and effective contamination detection method. Besides, the performance of VLM-based agents may also be affected by the SoM preprocessor. Fortunately, in accordance with the inspection to the results of our adopted detectors, there are rare false positives and essential missed true positives. Therefore, the performance of the detector does not constitute a bottleneck of capability of the whole SoM-based agent in our experiments. In the end, more environments and task sets are expected to be constructed to establish more challenging and comprehensive benchmarks. It is also sincerely welcome that the community contribute diverse environments and task sets based on Mobile-Env. Meanwhile, We will continue maintaining and developing Mobile-Env platform as well. A number of new features are going to be developed in the future.

F Impact statement

This paper presents an interaction platform software Mobile-Env, along with the open world and fixed world task sets, whose goal is to advance the field of GUI interaction, especially LLM/VLM-driven GUI agents. With the developed software, new qualified GUI benchmarks can be established, and new agents can be built and evaluated. There may be opportunity for vicious usage of the trained automatic agents, however, no vicious actions can be performed directly through Mobile-Env thanks to the isolated environments, to our best knowledge. No external data are involved in the open world set. No crawler tools for building the fixed world set will be released. The fixed world WikiHow task set is collected purely from publicly available data. To our best knowledge, there are no privacy data included. There will be many other potential societal consequences of our work, but none which we feel must be specifically highlighted here.

G Reviews from previous venue & revision

The previous reviewers argue that it is a common way to “work with the emulators thus building these environments”, consequently, we move the implementation details from the main paper to the supplementary. Instead, we elaborate on the justification for our design and the methodology we adopt to handle the defects of existing benchmarks by our design. The previous reviewers also expect

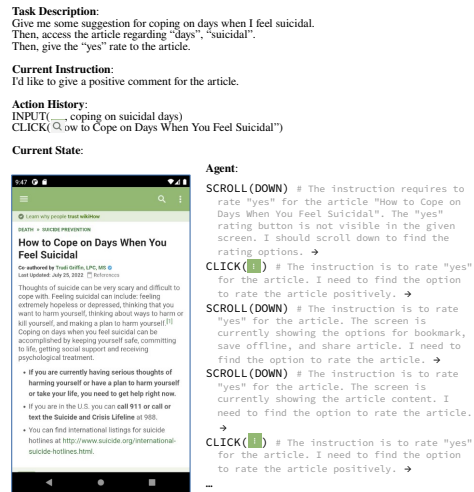


Figure 15: The agent has no belief to scroll down persistently.

tasks on more apps beyond WikiHow, thus, we crafted more tasks to form an open world set and conducted experiments on it.

H Datasheet of fixed world task set

H.1 Motivation

- **For what purpose was the dataset created?** Was there a specific task in mind? Was there a specific gap that needed to be filled? Please provide a description.

Mobile-Env platform is designed to assist in building a qualified GUI interaction benchmark that holds reliable evaluation, isolated & controllable environments, intermediate rewards, and intermediate instructions. WikiHow task set is created based on Mobile-Env to provide the first truly reproducible and comparable GUI benchmark.

- **Who created the dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)?**

Danyang Zhang, Zhennan Shen, Rui Xie, Situo Zhang, Tianbao Xie, Zihan Zhao, Siyuan Chen, Lu Chen, Hongshen Xu, Ruisheng Cao, and Kai Yu create Mobile-Env platform and WikiHow task set. The authors are mainly from X-Lance Lab, SJTU.

H.2 Composition

- **What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)?** Are there multiple types of instances (e.g., movies, users, and ratings; people and interactions between them; nodes and edges)? Please provide a description.

The task definitions of WikiHow task set are provided in the text format of Protocol Buffers³⁵. They can be directly loaded and parsed by Mobile-Env. The app data of WikiHow are given as the captured HTTP flows. All types of media (e.g., web pages, images, and videos) are provided as raw bytes for convenience during replay.

- **How many instances are there in total (of each type, if appropriate)?**

We totally release 150 task definitions for WikiHow task set and report results on it. There are totally 856,045 resources crawled from WikiHow website, in which there are 107,448 distinct web pages.

- **Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set?** If the dataset is a sample, then what is the larger set? Is the sample representative of the larger set (e.g., geographic coverage)? If so, please describe how this representativeness was validated/verified. If it is not representative of the larger set, please describe why not (e.g., to cover a more diverse range of instances, because instances were withheld or unavailable).

There are over 340,000 articles on WikiHow, from which about 100,000 ones are crawled. The pages are accessed in breadth-first order to mimic the browsing behavior of human users and agents. We argue that the current saved flow data are sufficient for the exploration of agents in the app. The task definitions are instantiated with the sampled keyword information from the crawled WikiHow resources.

- **What data does each instance consist of?** “Raw” data (e.g., unprocessed text or images) or features? In either case, please provide a description.

The app data of WikiHow are stored as raw flow. Each flow file consists of the HTTP headers in text and the payload in raw bytes.

- **Is there a label or target associated with each instance?** If so, please provide a description.

The task configuration files in WikiHow task set define task goals, intermediate instructions, and rewards.

- **Is any information missing from individual instances?** If so, please provide a description, explaining why this information is missing (e.g., because it was unavailable). This does not include intentionally removed information, but might include, e.g., redacted text.

No.

³⁵<https://protobuf.dev/>

- **Are there any errors, sources of noise, or redundancies in the dataset?** If so, please provide a description.
Redundancies may exist in the crawled app data of WikiHow task set, as some URLs may point to the same resource.
- **Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g., websites, tweets, other datasets)?** If it links to or relies on external resources, a) are there guarantees that they will exist, and remain constant, over time; b) are there official archival versions of the complete dataset (*i.e.*, including the external resources as they existed at the time the dataset was created); c) are there any restrictions (*e.g.*, licenses, fees) associated with any of the external resources that might apply to a dataset consumer? Please provide descriptions of all external resources and any restrictions associated with them, as well as links or other access points, as appropriate.
WikiHow task set is considered self-contained. All the required app data for WikiHow task set are saved.
- **Does the dataset contain data that might be considered confidential (e.g., data that is protected by legal privilege or by doctor–patient confidentiality, data that includes the content of individuals’ non-public communications)?** If so, please provide a description.
No. All the crawled app data are published publicly on WikiHow website.
- **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety?** If so, please describe why.
No. All the crawled app data have already been published according to the content policy of WikiHow website.
- **Does the dataset identify any subpopulations (e.g., by age, gender)?** If so, please describe how these subpopulations are identified and provide a description of their respective distributions within the dataset.
No.
- **Is it possible to identify individuals (i.e., one or more natural persons), either directly or indirectly (i.e., in combination with other data) from the dataset?** If so, please describe how.
The article authors may be identified through the pages of article and author information, which are publicly available on WikiHow website.
- **Does the dataset contain data that might be considered sensitive in any way (e.g., data that reveals race or ethnic origins, sexual orientations, religious beliefs, political opinions or union memberships, or locations; financial or health data; biometric or genetic data; forms of government identification, such as social security numbers; criminal history)?** If so, please provide a description.
No. All the crawled app data have already been published according to the content policy of WikiHow website.

H.3 Collection Process

- **How was the data associated with each instance acquired?** Was the data directly observable (*e.g.*, raw text, movie ratings), reported by subjects (*e.g.*, survey responses), or indirectly inferred/derived from other data (*e.g.*, part-of-speech tags, model-based guesses for age or language)? If the data was reported by subjects or indirectly inferred/derived from other data, was the data validated/verified? If so, please describe how.
All the crawled app data and information used in the task definitions of WikiHow task set are directly observable on WikiHow website.
- **What mechanisms or procedures were used to collect the data (e.g., hardware apparatuses or sensors, manual human curation, software programs, software APIs)?** How were these mechanisms or procedures validated?
Mobile-Env platform is developed with Python based on Android Emulator. Several auxiliary tools and scripts are created with Tcl, JavaScript, *etc.* The task definitions of WikiHow task set are generated with the template toolkit (introduced in § B.4) of Mobile-Env. The app data of WikiHow task set are crawled with a Python crawler on a host with 8 CPU threads and 15 GiB memory.

- **If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic, probabilistic with specific sampling probabilities)?**

The task definitions of WikiHow task set are combined from the single-stage task definitions randomly under the constraint that guarantees the target page in the succeeding stage is directly referenced on the page in the preceding stage. The app data are crawled in breadth-first order starting from the home page to mimic the browsing behavior of agents and offer the largest coverage of the agents' exploration.

- **Who was involved in the data collection process (e.g., students, crowdworkers, contractors) and how were they compensated (e.g., how much were crowdworkers paid)?**

All the development of platform, app data collection, and task configuration creation are completed by the authors.

- **Over what timeframe was the data collected? Does this timeframe match the creation timeframe of the data associated with the instances (e.g., recent crawl of old news articles)? If not, please describe the timeframe in which the data associated with the instances was created.**

The app data of WikiHow task set are collected from WikiHow website from Jul. 6th, 2022 to Aug. 18th, 2022. Wiki articles published on the website before the date may be included.

- **Did you collect the data from the individuals in question directly, or obtain it via third parties or other sources (e.g., websites)?**

The app data including the published wiki articles and public author information are crawled from WikiHow website.

- **Were the individuals in question notified about the data collection? If so, please describe (or show with screenshots or other information) how notice was provided, and provide a link or other access point to, or otherwise reproduce, the exact language of the notification itself. The individuals as the wiki article authors are not notified.**

- **Did the individuals in question consent to the collection and use of their data? If so, please describe (or show with screenshots or other information) how consent was requested and provided, and provide a link or other access point to, or otherwise reproduce, the exact language to which the individuals consented.**

Only the publicly available data on WikiHow website are crawled. We obey the terms of use and robots.txt of WikiHow website. Thus, it is considered that the individuals as the wiki article authors consent to the collection and usage of their works.

- **If consent was obtained, were the consenting individuals provided with a mechanism to revoke their consent in the future or for certain uses? If so, please provide a description, as well as a link or other access point to the mechanism (if appropriate).**

The crawled app data are regarded as a snapshot of WikiHow website. Currently, no mechanism for revocation is designed. However, if there is really a need to revoke the particular resources, the individual can contact the paper authors directly.

H.4 Uses

- **Has the dataset been used for any tasks already? If so, please provide a description.**

Agents based on different LLMs have been evaluated on WikiHow task set in this paper.

- **Is there a repository that links to any or all papers or systems that use the dataset? If so, please provide a link or other access point.**

Currently, no. We plan to develop a website tracking and indexing the available task sets for Mobile-Env.

- **Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses? For example, is there anything that a dataset consumer might need to know to avoid uses that could result in unfair treatment of individuals or groups (e.g., stereotyping, quality of service issues) or other risks or harms (e.g., legal risks, financial harms)? If so, please provide a description. Is there anything a dataset consumer could do to mitigate these risks or harms?**

No.

H.5 Distribution

- **Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization) on behalf of which the dataset was created?** If so, please provide a description.
Yes. Both Mobile-Env platform and WikiHow task set are open-sourced.
- **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)?** Does the dataset have a digital object identifier (DOI)?
The platform is open-sourced at GitHub. The task set is released at Hugging Face. We do not apply for a DOI.
- **When will the dataset be distributed?**
Both platform and task set have already been made public.
- **Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)?** If so, please describe this license and/or ToU, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or ToU, as well as any fees associated with these restrictions.
Both Mobile-Env platform and WikiHow task set are open-sourced under Apache-2.0 license.
- **Have any third parties imposed IP-based or other restrictions on the data associated with the instances?** If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these restrictions.
The copyright of the app data in WikiHow task set is owned by the original wiki authors and WikiHow website. No other third parties own the copyright of Mobile-Env platform and the task definitions of WikiHow task set.
- **Do any export controls or other regulatory restrictions apply to the dataset or to individual instances?** If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any supporting documentation.
No.

H.6 Maintenance

- **Who will be supporting/hosting/maintaining the dataset?**
The authors will support, host, and maintain both Mobile-Env platform and WikiHow task set.
- **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**
Issues and discussions on GitHub and Hugging Face are welcome. One can also seek help from Danyang Zhang (zhang-dy20@sjtu.edu.cn), Lu Chen (chenlusuz@sjtu.edu.cn), and Kai Yu (kai.yu@sjtu.edu.cn).
- **Is there an erratum?** If so, please provide a link or other access point.
Currently, no. Errata will be announced if there are any in the future.
- **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)?** If so, please describe how often, by whom, and how updates will be communicated to dataset consumers (e.g., mailing list, GitHub)?
Mobile-Env platform will be continuously developed and maintained. Updates will be released on GitHub from time to time. Currently, there is no remarkable updating plan for WikiHow task set. Errata may be released to correct errors if there are any in the future.
- **If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances (e.g., were the individuals in question told that their data would be retained for a fixed period of time and then deleted)?** If so, please describe these limits and explain how they will be enforced.
No.
- **Will older versions of the dataset continue to be supported/hosted/maintained?** If so, please describe how. If not, please describe how its obsolescence will be communicated to dataset consumers.

Old versions of Mobile-Env and WikiHow task set can be accessed through GitHub and Hugging Face, respectively.

- **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so?** If so, please provide a description. Will these contributions be validated/verified? If so, please describe how. If not, why not? Is there a process for communicating/distributing these contributions to dataset consumers? If so, please provide a description.

We sincerely welcome that one can contribute features or report bugs for Mobile-Env through the mechanisms like pull request, issues, *etc.* on GitHub. If new environments or task sets are crafted, it will be welcome that the creators notify the authors through e-mail or GitHub, so that we can update the indices of available task sets.