

---

# GUI-Robust: A Comprehensive Dataset for Testing GUI Agent Robustness in Real-World Anomalies

---

**Jingqi Yang\***  
Zhejiang University

**Zhilong Song**  
Intelligence Indeed

**Jiawei Chen**  
Zhejiang University

**Mingli Song**  
Zhejiang University

**Sheng Zhou**  
Zhejiang University

**linjun sun**  
Intelligence Indeed

**Xiaogang Ouyang**  
Intelligence Indeed

**Chun Chen**  
Zhejiang University

**Can Wang**  
Zhejiang University

## Abstract

The development of high-quality datasets is crucial for benchmarking and advancing research in Graphical User Interface (GUI) agents. Despite their importance, existing datasets are often constructed under idealized conditions, overlooking the diverse anomalies frequently encountered in real-world deployments. To address this limitation, we introduce GUI-Robust, a novel dataset designed for comprehensive GUI agent evaluation, explicitly incorporating seven common types of anomalies observed in everyday GUI interactions. Furthermore, we propose a semi-automated dataset construction paradigm that collects user action sequences from natural interactions via RPA tools and then generate corresponding step and task descriptions for these actions with the assistance of MLLMs. This paradigm significantly reduces annotation time cost by a factor of over 19 times. Finally, we assess state-of-the-art GUI agents using the GUI-Robust dataset, revealing their substantial performance degradation in abnormal scenarios. We anticipate that our work will highlight the importance of robustness in GUI agents and inspires more future research in this direction. The dataset and code are available at <https://github.com/chessbean1/GUI-Robust>.

## 1 Introduction

A GUI agent is an intelligent system capable of autonomously interacting with graphical user interfaces by perceiving visual elements, understanding task objectives, and executing corresponding actions [15]. The development of GUI agents holds significant promise for automating GUI operations, reducing human workload and enabling seamless human-computer collaboration. Traditionally, GUI automation has relied on heuristic rules, brittle scripts, or hard-coded templates, which limit flexibility and generalization [9] [20]. Recently, the emergence of Multimodal Large Language Models (MLLMs) [2] [16] [22] has accelerated GUI agent research. By integrating visual perception with language understanding, MLLMs enable agents to interpret complex screen layouts, comprehend natural language instructions, and reason about sequential actions in open-ended environments [6] — marking a significant advance toward fully GUI automation.

In this rapidly evolving field, the construction of high-quality datasets is essential for benchmarking and advancing GUI agent research. However, existing datasets are typically constructed under

---

\*Contact: chessbean1@gmail.com

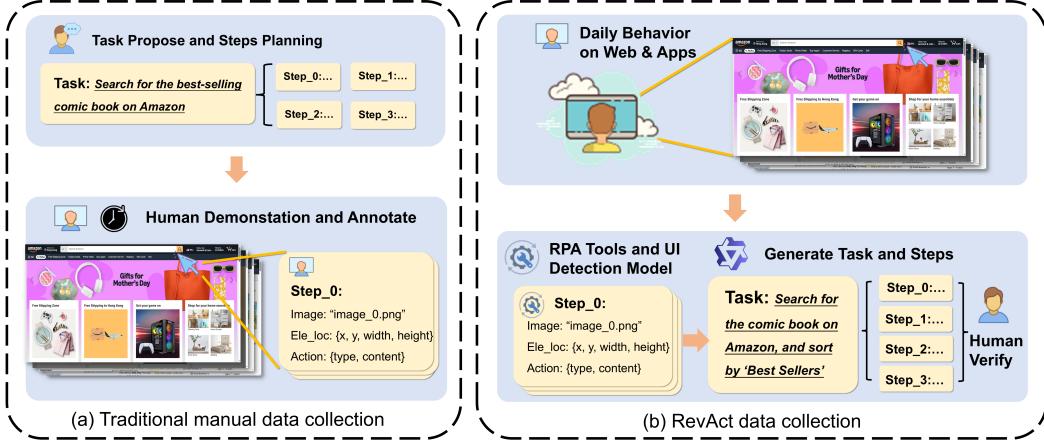


Figure 1: Comparison between (a) traditional manual data collection and (b) our proposed RevAct pipeline. Manual pipelines require annotators to explicitly design tasks, plan execution steps, demonstrate actions, and annotate each step. In contrast, RevAct records natural user behavior on web/apps, leverages automated tools (e.g., RPA, detection models, and LLMs) to infer action semantics and generate step/task descriptions, at last verify by human.

idealized conditions and fail to capture the diverse range of anomalies encountered in real-world deployment. In practice, GUI agents deployed in industrial or consumer applications frequently encounter unexpected failure modes and environmental disturbances, such as action failures, obstructive pop-up advertisements, network disconnections, etc. These anomalies can significantly disrupt execution flow, leading to task failures or, in severe cases, unintended interactions with sensitive interface components, potentially resulting in erroneous or hazardous outcomes.

**New Dataset: GUI-Robust.** To address this gap, we present **GUI-Robust**, the first dataset designed to evaluate robustness in the presence of abnormal scenarios. GUI-Robust contains 5,318 annotated tasks (consists of task descriptions and user behaviors sequences) collected from 392 diverse sources, spanning both websites and third-party desktop applications on Windows. Notably, it includes 200 abnormal tasks covering 7 common types of anomalies encountered in everyday GUI usage including action failure, login page, captcha page, ad pop-up, cookie pop-up, page loading and network disconnection. This enables rigorous evaluation of GUI agent robustness against real-world anomalies, a critical aspect for practical deployment.

Beyond its unique property of robustness, as shown in Table 1, GUI-Robust offers several advantages for comprehensive evaluation: (1) it includes a broad range of task and action types including click, input text, retrieve information from page, open a new web or app and report anomalies to human; (2) it incorporates cross-scenario tasks spanning multiple applications or websites, reflecting more realistic and complex workflows; and (3) it covers both Chinese and English software environments. A comparison between GUI-Robust and existing benchmarks is shown in Table 1. These features ensure that our dataset closely aligns with real-world usage, supporting thorough assessment of agent robustness, adaptability, and generalization.

**Novel Data Collection Strategy: RevAct.** Traditional dataset construction typically relies on a manual pipeline — ranging from task design to execution planning and manual demonstration — which is labor-intensive, costly, and requires significant domain expertise. To address these limitations, we propose a novel data collection strategy, RevAct, which reverses the conventional workflow and enables semi-automated dataset generation: we first collect user action sequences from natural interactions via RPA(Robotic Process Automation) tools, and then generate specific step and task descriptions for these actions with the assistance of MLLMs. This approach substantially reduces annotation costs, as expert involvement is limited to reviewing and revising step and task descriptions generated by MLLMs. Specifically, for each user action sequence, we leverage YOLOv8 for GUI element detection and Qwen2.5-VL for task generation and summarization, achieving over 71% accuracy in automatic task generation. Human annotators are only required for minimal correction, reducing annotation time by a factor of over 19 times compared to the traditional pipeline.

Table 1: Statistics of GUI-Robust compared with existing datasets.

Benchmark	#Apps/Web	#Tasks	Platform	Abnormal?	Cross?	Lang.
Mind2Web	137	2,350	Desktop Web	✗	✗	EN
WebLINX	155	2,337	Desktop Web	✗	✗	EN
OSWorld	9	369	Desktop(Apps+Web)	✗	✓	EN
WindowsAgentArena	11	154	Desktop(Apps+Web)	✗	✗	EN
GUI Odyssey	201	7,735	Android Apps	✗	✓	EN
SPA-Bench	58	340	Android Apps	✗	✓	EN & CN
WorldGUI	10	315	Desktop Apps	✗	✗	EN
GUI-Robust	392	5,318	Desktop(Apps+Web)	✓	✓	EN & CN

**Comprehensive Experiments.** GUI-Robust provides a thorough benchmark for evaluating GUI agents on element grounding, multi-step task completion, cross-scenario execution, and robustness under abnormal conditions. We conduct extensive experiments with three representative MLLMs and two specific-designed GUI agents, demonstrating that all models experience significant performance degradation in abnormal scenarios. This underscores that robust GUI understanding in real-world settings remains a challenging problem. We hope our work will draw greater attention to the robustness of GUI agents and inspire further research in this direction based on our dataset.

## 2 Ralated Work

### 2.1 GUI Agents

Early progress in GUI agent research has been largely driven by the rapid advancement of multimodal large language models (MLLMs), such as GPT-4 [17], Gemini-2.5 [22], Qwen-VL [1] and Claude3.5. These foundation models—capable of processing both visual and textual inputs—enable agents to interpret screen layouts, follow natural language instructions, and reason about GUI actions in open-ended environments. Recently, there has been a surge in the development of GUI-specific agents, which fine-tune MLLMs or build on their architecture to improve performance in structured GUI environments. Notable examples include: CogAgent [10], built on the GLM-VLM [8] framework, which achieves improvements in perception, action space coverage, and reasoning via staged training; ShowUI [12], trained on top of Qwen2-VL [23], with enhancements for UI tokenization and cross-system navigation; UI-TARS [18], which integrates Qwen2-VL [23] with additional task planning and grounding modules, excelling in complex GUI manipulation tasks. In particular, two fundamental capabilities have emerged as essential for effective GUI automation: element grounding, which refers to the ability to accurately locate relevant UI elements on the screen, and task completion, which involves determining the correct type and content of interaction at each step.

### 2.2 GUI Benchmarks

Benchmark datasets are fundamental to the development and evaluation of GUI agents. A key capability required for effective GUI interaction is accurate element grounding. For this purpose, ScreenSpot[5] serves as a dedicated benchmark, offering both mobile and web-based examples to evaluate agents’ spatial grounding capabilities. For task completion, Mind2Web[7] stands out as a pioneering benchmark built in real-world web environments rather than simplified simulations. WebVNLN[4] focuses on navigation and QA within shopping websites, while WebLINX[14] introduces multi-turn interactions in realistic web interfaces. Recently, some datasets have emerged that move toward dynamic GUI environment. OSWorld[25] and WindowsAgentArena[3] are two benchmark platforms targeting GUI agents in Ubuntu and Windows operating systems, respectively. However, these dynamic testing are time-consuming and demand high setup Cost. In domain-specific research, WebWalkerQA[24] emphasizes web-based information retrieval and question answering. On the Android platform, PIXELHELP[11] and AITW (Android in the Wild)[19] offer static mobile task datasets, while META-GUI[21] focuses on multi-turn dialog interactions. GUI Odyssey[13] proposes cross-app task scenarios. Despite these advances, existing datasets are mostly constructed under ideal conditions, failing to address the irregularities and exceptions commonly found in real-world industrial scenarios. These abnormal situations are critical for evaluating agent **robustness**. WorldGUI[26] attempts to address this by exploring the impact of varying initial states, but many real-world

Table 2: Composition of the GUI-Robust Dataset

Category	Description	#Tasks	Action Space
Standard Tasks	routine tasks with single web/app	4414	click, input
Information Retrieval Tasks	tasks with information retrieval	503	click, input, get_info
Cross-Scenario Tasks	tasks across two webs/apps	201	click, input, get_info, open
Abnormal Tasks	tasks with abnormal scenarios	200	click, input, get_info, wait, human

anomalies remain unaddressed. To fill this gap, we introduce **GUI-Robust**, a benchmark specifically designed to evaluate agent robustness under abnormal GUI scenarios.

### 2.3 Data Collection Method for GUI Benchmark

Most existing GUI datasets are collected with a manual pipeline ranging from task design to execution planning and manual demonstration. For instance, Mind2Web[7] outlines a process involving *Website Selection*, *Task Proposal*, *Task Demonstration*, and *Task Verification*, all of which heavily depend on human effort. Similarly, WebLINX[14] requires annotators to manually design tasks and engage in complex dialogues or task edits. To improve efficiency, some datasets leverage LLMs to accelerate task designs. Mind2Web[7] uses ChatGPT to generate seed tasks that inspire human annotators. GUI Odyssey[13] employs a template-based approach, where annotators design reusable templates and GPT-4 replaces entities (*e.g.*, app or item names) to generate multiple variants. Remarkably, these methods are still labor-extensive as it only accelerates task design procedure and still requires annotors for actions planning and execution, which is time-consuming. A more autonomous approach is seen in WebWalkerQA[24], where GPT-4o[16] generates QA pairs based solely on webpage screenshots. While promising, this method is currently limited to QA-type tasks and does not generalize to broader GUI tasks. To overcome these limitations, we propose a reverse data collection paradigm that collects user action sequences from natural interactions and then generate corresponding step and task descriptions for these actions with the assistance of MLLMs. Human annotators are only required for reviewing and revising step and task descriptions, which is highly efficient.

## 3 GUI-Robust Dataset

### 3.1 Overview

GUI-Robust comprises 5,318 annotated tasks collected from 392 distinct websites and desktop applications. Each task (*i.e.*, data instance) consists of a task description and a sequence of user behaviors. Notably, GUI-Robust includes 200 tasks under abnormal scenarios, encompassing seven distinct types of anomalies commonly encountered in everyday GUI interactions. This enables rigorous evaluation of GUI agent robustness under real-world conditions. In addition to this unique feature, GUI-Robust offers a greater diversity and complexity of tasks, such as those involving information retrieval and cross-scenario operations. Furthermore, our dataset spans a wide range of 49 domains and covers nearly all commonly used platforms within the Chinese internet ecosystem. These characteristics ensure that GUI-Robust is closely aligned with real-world usage, thereby supporting comprehensive assessment of agent robustness, adaptability, and generalization. The dataset statistics are presented in Table 1 and 2.

To promote ease of use, we release a lightweight evaluation toolkit for GUI-Robust (see details in Appendix D), which enables users to quickly evaluate existing models or integrate their own with minimal effort.

### 3.2 Metadata Definition

**Task Formulation.** For a specific task, the execution of by a GUI agent can be represented as a sequence of actions paired with corresponding UI element coordinates at each step, as an example illustrated in Fig. 3. Formally, given a task description  $T$  and a screenshot  $S_i$  at step  $i$ , the agent  $\mathcal{G}$  generates an action  $A_i$  and selects a UI element  $E_i = (x_i, y_i)$  located at screen coordinates  $(x, y)$ . This step-wise interaction is defined as:

$$R_i = (A_i, E_i) = \mathcal{G}(S_i, T)$$

Upon task completion, the full execution result is represented as a sequence:

$$R = \{(A_i, E_i) |_{i=1}^n, T\}$$

where  $n$  denotes the total number of steps required to complete the task.

**Dataset Structure.** GUI-Robust is composed of a collection of annotated tasks, each comprising a task description  $T$  and a sequence of screenshots  $S = \{S_1, S_2, \dots, S_n\}$ . The task description  $T$  is provided in natural language, reflecting realistic user instructions (e.g., “Search for the best-selling comic book on Amazon”), while the sequence  $S$  represents the screenshots captured at each interaction step. For each screenshot  $S_i$ , we also provide the corresponding ground-truth step description (e.g., “Click on the ‘Account & Lists’ text”), the action performed (e.g., “Click”), and the element location. These annotations facilitate the evaluation of the agent’s element grounding and task completion capabilities. See Appendix A for details on the data format structure.

**Element Type** In this dataset, we categorize UI elements into three types: **icon**, **text**, and **box**. The icon type represents interactive UI elements that are abstract icons with no textual content. The text type refers to UI elements containing text, such as buttons or links. The box type represents input fields, where the task involves entering content into a text input box. This classification allows us to assess the localization capabilities of different models across various types of UI elements, providing insights into their performance in identifying and interacting with different element categories.

**Action Space** The action space in **GUI-Robust** comprises six distinct action types: `click`, `input`, `get_info`, `open`, `wait`, and `human`. Among these, `click`, `input`, `get_info` and `open` are four standard actions, corresponding to “click an element”, “input the content”, “retrieve information from the interface”, “open a the website or application”. In addition to these standard actions, GUI-Robust incorporates two specialized action types to address anomalous scenarios: `wait` denotes the “wait for response” and `human` denotes the “require human intervention”. the more details about these actions refer to Appendix A.

### 3.3 Abnormal Scenarios

By collecting and analyzing 5,318 real user interaction sessions across a diverse range of websites and desktop applications, we identified seven types of commonly encountered anomalies (See Fig. 4):

- **Action Failure:** The agent’s previous action does not trigger the expected UI response (e.g., a button click has no effect).
- **Login Page:** A login prompt appears unexpectedly, requiring authentication before proceeding.
- **Captcha Page:** The agent encounters a CAPTCHA challenge that it cannot autonomously solve.
- **Ad Pop-up:** Advertisements pop up and obscure key UI elements, disrupting the intended interaction flow.
- **Cookie Consent Pop-up:** A cookie consent dialog appears and must be dismissed before interacting with the main interface.
- **Page Loading Delay:** The page remains in a loading state for an extended time, preventing access to target elements.
- **Network Disconnection:** The interface fails to load due to temporary or complete loss of network connectivity.

These scenarios emerged organically from natural human activity on the internet and reflect common disruptions that hinder successful task completion. We introduce `wait` and `human` actions for agents to solve or propose the problem (See Table 6). This observation motivated us to explicitly incorporate these real-world failure cases into our dataset design, in order to enable more robust benchmarking and guide the development of agents that can operate reliably in practical environments.

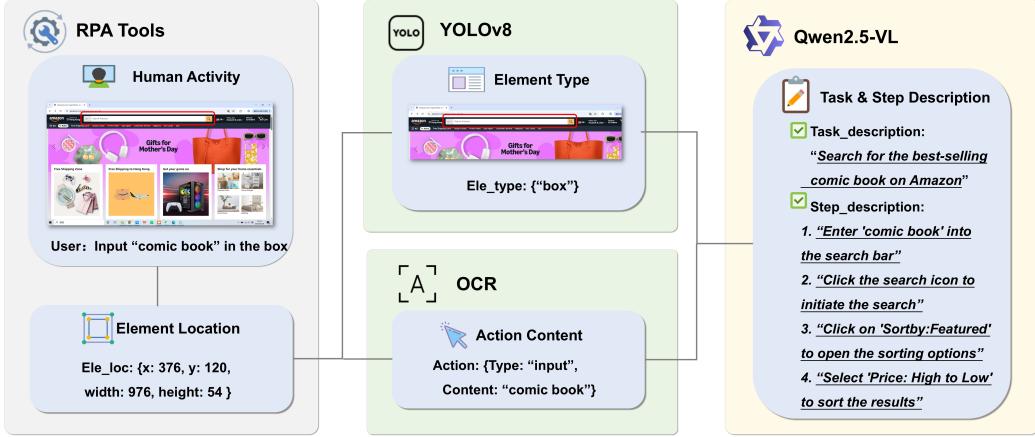


Figure 2: Overview of the RevAct pipeline for semi-automated data collection. Human interaction traces are captured by RPA tools, then processed by a YOLOv8 detector and OCR to extract element type, location, and action content. A multimodal LLM (Qwen2.5-VL) generates step-wise and task-level descriptions, enabling efficient and scalable construction of annotated GUI tasks.

## 4 Data Collection Method

### 4.1 Existing Limitations

Both step and task descriptions and corresponding action sequences are recorded to construct the dataset. However, this approach is labor-intensive and inefficient, particularly in the following two respects: 1) Task design typically depends on human heuristics, making the process cognitively demanding and time-consuming. Furthermore, the scope of manually designed tasks is often limited, and the curated task patterns frequently diverge from real-world user behaviors. 2) For recording action sequences, human annotators are required to have substantial familiarity with the target websites and desktop applications, which necessitates specialized training and often multiple iterations to produce accurate records. Even for annotators with considerable expertise, our experiments demonstrate that the average annotation time per task exceeds 15 minutes.

### 4.2 New Data Collection Strategy: RevAct

To address these challenges, we propose a new data collection pipeline that reverses the convention workflow and enables semi-automated dataset generation. The basic idea is to collect user action sequences from natural interactions, and then generate specific step and task descriptions for these actions with the assistance of MLLMs. This strategy could avoid the time-consuming steps of the task design and manual demonstration, requiring Human annotators only involves to review and revise step and task descriptions generated by MLLMs. The pipeline is illustrated in Fig. 2. Specifically, the data collection consists of the following three steps:

- 1. Capture — Screenshot and Element Coordinate Capture:** We begin by collecting user action sequences derived from natural interactions. Specifically, we recruit volunteers and record their routine activities on familiar websites and desktop applications. During these sessions, a Robotic Process Automation tool (Indeed-Intelligence RPA<sup>2</sup>) automatically captures both screenshots and the coordinates of UI elements for each user action. In contrast to traditional approaches, which require annotators to perform pre-defined tasks demanding expertise and extensive planning, our data collection process is high efficient and does not necessitate specialized knowledge. This step can be easily and efficiently completed for various volunteers without expertise, does not require expertise annotators. The economic cost associated with this procedure is negligible. Besides, this method ensures that the collected data more accurately reflects real-world user behavior.
- 2. Interpretation — Action Recognition:** The recorded screenshots and element coordinates are fed into a YOLOv8-based GUI element detection model, which identifies the type of UI element

<sup>2</sup><https://www.ai-indeed.com/>

interacted with at each step. Subsequently, we establish a one-to-one mapping between element types and action types. For example:

- If the operated UI element is a clickable icon or text, the action type is mapped to `click`.
- If the UI element is an input box, the action type is mapped to `text input`.

Unfortunately, we observed that YOLOv8 demonstrates limited accuracy in recognizing other types of actions. However, the number of other action types is relatively small, and we simply leave them for subsequent revision. Finally, we use an OCR tool to extract the content of the UI element, supplementing the action’s semantic information.

3. **Summarize — Step and Task Description Generation:** We feed the previously collected screenshots and the identified actions into Qwen2.5-VL [2]. The model generates a step-by-step description based on this information, then summarizes all step descriptions to produce a task-level description, thereby completing the construction of a fully annotated data instance. In the prompt provided to the multimodal model, we explicitly instruct it to follow two key principles: (1) The task description should be a natural language summary of the overall task goal, rather than a simple concatenation of the step descriptions. (2) The description should emphasize the outcome of the final step, as this often captures the core objective of the task. At the same time, it should incorporate relevant details from intermediate steps to ensure that the resulting instruction is both semantically accurate and operationally complete.
4. **Revise — Refine the Step and Task Descriptions :** We employ human annotators to review and refine the step and task descriptions generated by MLLMs. Given that MLLMs exhibit relatively high accuracy, this process is efficient, as annotators are required to revise only a small subset of the descriptions.

### 4.3 Empirical Results

To validate the effectiveness of the proposed method, we conducted a feasibility evaluation by having human annotators review and correct the task data generated by the RevAct method. The evaluation process focused on the following three aspects:

- **Accuracy of step descriptions:** Whether each step description accurately reflects the UI element and its content corresponding to the operation.
- **Consistency of task descriptions:** Whether the task description aligns with the actual intent and all steps.
- **Completeness of the task description:** Whether the task description contains sufficient information to reproduce the full interaction sequence—i.e., whether a human can follow it to reproduce the entire sequence of actions.

In the experiment, we test the descriptions generated by the RevAct method and invited 10 experienced annotators to independently correct and assess the quality. The results showed that, for standard tasks in Chinese, the step and task description accuracy was 78.63% and 72.57%. In comparison, for standard tasks in English, the step and task description accuracy reached 71.71% and 86.67%. These results indicate that the RevAct method can generate accurate operation and task text across different languages and task complexities.

In terms of data collection efficiency, RevAct significantly reduced the workload of annotators compared to traditional fully manual methods. Specifically, collecting and processing 100 task samples using traditional methods (See details in Appendix E) typically requires 10 annotators and takes about 150 minutes, whereas the RevAct method, required only 10 annotators and 7.8 minutes to review and finalize the same number of collected samples. This represents an over  $19\times$  improvement in efficiency, highlighting the method’s significant advantage in reducing human labor costs and accelerating data construction.

## 5 Experiment

### 5.1 Experiment Setup

We evaluate 5 representative models on the **GUI-Robust** dataset, including 3 general-purpose multimodal large models with visual input and understanding capabilities—GPT4o [16], Qwen2.5-

Table 3: Performance on the UI Element Grounding. **Bold** indicates the best performance across all models, while underlined values highlight the best among general-purpose multimodal models (MLLMs).

Model	Action Acc.	Ele_loc Acc.			
		Icon_ele	Text_ele	Box_ele	All
Qwen2.5-vl-7B	82.59	4.63	<u>5.07</u>	4.13	4.77
Qwen2.5-vl-72B	89.20	<u>5.56</u>	3.21	1.83	3.40
Gemeini2.5-Flash-Preview	77.24	1.85	3.72	<u>22.02</u>	<u>7.20</u>
GPT-4o-2024-11-20	<b>90.56</b>	1.39	3.04	5.50	3.21
UI-TARS-1.5-7B	51.17	76.85	<b>84.63</b>	62.39	<b>78.31</b>
CogAgent-9B	42.40	<b>78.65</b>	67.67	<b>89.50</b>	77.12

VL [2] and Gemini 2.5-Flash-Preview [22]—and two pretrained agents explicitly designed for GUI automation—UI-TARS [18] and CogAgent [10].

We design 2 evaluation types, comprising a total of 5 evaluation tasks:

**UI Element Grounding:** We randomly sample 1,028 single-step trajectories from the *Standard Tasks* subset. Given a step description and its corresponding screenshot, models are evaluated on: (a) **Action Accuracy (Action Acc.)**: both action type and content must match the ground truth. (b) **Coordinate Accuracy (Coord. Acc.)**: predicted coordinates must fall within the bounding box of the target UI element.

**Task Completion:** We assigned 5 models to execute a set of tasks sampled from GUI-Robust: **302 Standard Tasks** (randomly selected from the 4,925-task subset) and **197 Abnormal Tasks**. During Task Completion evaluation, the model receives three inputs at each step: (a) the task description, (b) the current page screenshot, and (c) a history of previously predicted actions and UI element coordinates. Based on this context, the model is required to generate the next action and its corresponding element location, progressing until the entire task is completed. We evaluate: (a) Action Acc. (2) Coord. Acc. and (c) **Task Success Rate (SR)**: the task is considered successful only if all actions and coordinates match the ground truth across the full trajectory.

## 5.2 Results

**UI Element Grounding.** Table 3 presents the results of various models on the UI Element Grounding task. We draw several key observations:

1. UI-TARS and CogAgent significantly outperform general-purpose MLLMs on grounding accuracy. UI-TARS achieves the highest overall localization accuracy (78.31%) across all element types, followed closely by CogAgent (77.12%). This demonstrates the advantage of GUI-specific pretrained models in accurately grounding UI elements, especially for visually complex or fine-grained components. All MLLMs underperform significantly in overall localization ( $\leq 7.20\%$ ), highlighting a persistent gap in visual grounding capabilities when applied to GUI domains;
2. General-purpose MLLMs like GPT-4o and Qwen2.5 show strong action recognition but weak spatial grounding. GPT-4o-2024 achieves the best action accuracy at 90.56%, and Qwen2.5-vl-72B follows with 89.20%. However, both models struggle to localize elements precisely—e.g., GPT-4o achieves only 3.21% overall localization accuracy, suggesting that while it understands what to do, it often fails to locate where to do it;
3. Different models specialize in different UI element types. CogAgent excels at localizing input boxes (89.50%) and icon elements (78.65%). UI-TARS performs best on text-based elements (84.63%). Interestingly, Gemini2.5-Flash achieves the highest accuracy on box elements (22.02%) among the MLLMs, despite lower overall performance.

**Task Completion.** Table 4 presents a comparative analysis of model performance on full-task execution under both normal and abnormal conditions. Several key observations emerge:

1. Under normal conditions, GUI-specific agents outperform general-purpose models. UI-TARS and CogAgent demonstrate significantly higher task success rates (21.61% and 14.38%, respectively)

Table 4: Performance on Task Completion evaluation. **Bold** indicates the best performance across all models, while underlined values highlight the best among general-purpose multimodal models (MLLMs).

Model	Normal			Abnormal		
	Action Acc.	Ele_loc Acc.	Task SR	Action Acc.	Ele_loc Acc.	Task SR
Qwen2.5-VL-72B	<b>56.99</b>	3.03	<u>1.32</u>	48.54	0.98	0
Gemini2.5-flash-preview	50.63	<u>6.78</u>	0	45.15	<u>5.13</u>	0
GPT-4o-2024-11-20	56.26	1.88	0.33	<b>55.07</b>	1.95	0
UI-TARS-1.5-7B	35.36	<b>62.33</b>	<b>21.61</b>	28.88	<b>61.77</b>	<b>4.73</b>
CogAgent-9B	33.12	54.32	14.38	23.48	44.35	2.14

and element localization accuracy (62.33% and 54.32%) than LLM-based models such as GPT-4 or Qwen2.5-VL. This suggests that specialized training on GUI interactions yields clear benefits in structured, expected environments.

2. In abnormal settings, general-purpose LLMs show stronger anomaly awareness, but poor spatial grounding. Despite lower task success rates, models like GPT-4o and Qwen2.5-VL maintain higher action accuracy (55.07% and 48.54%, respectively) compared to GUI-specific agents. This indicates that LLMs are more sensitive to abnormal cues (e.g., detecting login pop-ups or page errors), but fail to locate or interact with the appropriate UI elements—limiting their ability to recover or proceed correctly.

GUI-specific agents struggle with unanticipated anomalies. These models often lack mechanisms to recognize or adapt to unexpected interface changes. For example, during a task on the “58.com” website, UI-TARS encounters a sudden login page. Instead of detecting and reporting the abnormal transition, it predicts the action “click 58.com”, which will cause an infinite redirection loop. This failure illustrates the rigidity of GUI-specific models when exposed to disruptions outside their training distribution.

Incorrect behavior under anomalies can lead to critical consequences. In real-world applications, failure to handle abnormal states may result in hazardous operations such as accidental deletions, repeated entry of sensitive credentials, or navigation to malicious or unintended pages. Such behaviors not only reduce task reliability but also raise serious safety and data integrity concerns. These results highlight the need for robust anomaly-handling capabilities in GUI agents—not only to succeed in idealized environments, but also to detect, report, and gracefully recover from irregularities.

## 6 Discussion

**Limitations.** While GUI-Robust offers a substantial step toward evaluating agent robustness in realistic GUI environments, several limitations remain. The abnormal scenarios in our dataset—though diverse—are still limited to 7 predefined categories. Real-world applications may involve more complex or compound failure modes (e.g., cascading authentication prompts, adaptive overlays, operations with insufficient permissions), which are not currently covered. Second, GUI-Robust adopts a static dataset and evaluation paradigm, where each task is annotated with a single canonical execution trajectory.

**Future Work.** Future extensions of GUI-Robust aim to address these limitations along two directions. First, we plan to enrich the spectrum of abnormal scenarios by including more diverse failure types—such as dynamic Captcha variants, pop-ups requiring multi-step dismissal, and access control violations. Second, we envision evolving GUI-Robust into a dynamic evaluation platform, where agents interact with simulated environments rather than static trajectories. In such settings, evaluation criteria can move beyond exact action or coordinate matching to focus on goal-oriented success, state transition correctness, and error recovery ability. This shift will better reflect real-world deployment conditions and enable the development of agents that are not only accurate but also flexible, adaptive, and robust.

## 7 Conclusion

In this work, we present GUI-Robust, a comprehensive benchmark designed to evaluate the robustness of GUI agents under real-world anomalies, which incorporates 7 types of common abnormal scenarios. We further introduce RevAct, an efficient upside-down data collection pipeline that leverages multimodal models to generate rich, naturalistic task annotations from user behavior traces. Through extensive experiments we demonstrate that existing GUI-specific agents struggle in failure-prone environments. These findings highlight the pressing need for robustness-focused datasets and evaluation protocols. We hope GUI-Robust will serve as a valuable resource for advancing the development of reliable, adaptable, and safe GUI agents in real-world applications.

## References

- [1] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023.
- [2] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- [3] Rogerio Bonatti, Dan Zhao, Francesco Bonacci, Dillon Dupont, Sara Abdali, Yinheng Li, Justin Wagle, Kazuhito Koishida, Arthur Bucker, Lawrence Jang, and Zack Hui. Windows agent arena: Evaluating multi-modal os agents at scale. September 2024.
- [4] Qi Chen, Dileepa Pitawela, Chongyang Zhao, Gengze Zhou, Hsiang-Ting Chen, and Qi Wu. Webvln: Vision-and-language navigation on websites. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2024.
- [5] Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Li YanTao, Jianbing Zhang, and Zhiyong Wu. SeeClick: Harnessing GUI grounding for advanced visual GUI agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9313–9332, Bangkok, Thailand, August 2024. Association for Computational Linguistics. URL <https://aclanthology.org/2024.acl-long.505>.
- [6] Daniel Chin, Yuxuan Wang, and Gus Xia. Human-centered llm-agent user interface: A position paper, 2024. URL <https://arxiv.org/abs/2405.13050>.
- [7] 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, 2023.
- [8] Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie Tang, Jing Zhang, Juanzi Li, Lei Zhao, Lindong Wu, Lucen Zhong, Mingdao Liu, Minlie Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shudan Zhang, Shulin Cao, Shuxun Yang, Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan Zhang, Xiaotao Gu, Xin Lv, Xinghan Liu, Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan An, Yifan Xu, Yilin Niu, Yuantao Yang, Yueyan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu Wang, Zhen Yang, Zhengxiao Du, Zhenyu Hou, and Zihan Wang. Chatglm: A family of large language models from glm-130b to glm-4 all tools, 2024.
- [9] Theodore Hellmann and Frank Maurer. Rule-based exploratory testing of graphical user interfaces. pages 107 – 116, 09 2011. doi: 10.1109/AGILE.2011.23.
- [10] Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxiao Dong, Ming Ding, and Jie Tang. Cogagent: A visual language model for gui agents, 2023.

- [11] Yang Li, Jiacong He, Xin Zhou, Yuan Zhang, and Jason Baldridge. Mapping natural language instructions to mobile ui action sequences. In *Annual Conference of the Association for Computational Linguistics (ACL 2020)*, 2020. URL <https://www.aclweb.org/anthology/2020.acl-main.729.pdf>.
- [12] Kevin Qinghong Lin, Linjie Li, Difei Gao, Zhengyuan Yang, Shiwei Wu, Zechen Bai, Weixian Lei, Lijuan Wang, and Mike Zheng Shou. Showui: One vision-language-action model for gui visual agent, 2024. URL <https://arxiv.org/abs/2411.17465>.
- [13] Quanfeng Lu, Wenqi Shao, Zitao Liu, Fanqing Meng, Boxuan Li, Botong Chen, Siyuan Huang, Kaipeng Zhang, Yu Qiao, and Ping Luo. Gui odyssey: A comprehensive dataset for cross-app gui navigation on mobile devices. *arXiv preprint arXiv:2406.08451*, 2024.
- [14] Xing Han Lù, Zdeněk Kasner, and Siva Reddy. Weblinx: Real-world website navigation with multi-turn dialogue, 2024.
- [15] Dang Nguyen, Jian Chen, Yu Wang, Gang Wu, Namyong Park, Zhengmian Hu, Hanjia Lyu, Junda Wu, Ryan Aponte, Yu Xia, Xintong Li, Jing Shi, Hongjie Chen, Viet Dac Lai, Zhouhang Xie, Sungchul Kim, Ruiyi Zhang, Tong Yu, Mehrab Tanjim, Nesreen K. Ahmed, Puneet Mathur, Seunghyun Yoon, Lina Yao, Branislav Kveton, Thien Huu Nguyen, Trung Bui, Tianyi Zhou, Ryan A. Rossi, and Franck Dernoncourt. Gui agents: A survey, 2024. URL <https://arxiv.org/abs/2412.13501>.
- [16] OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Mądry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex Renzin, Alex Tachard Passos, Alexander Kirillov, Alexi Christakis, Alexis Conneau, Ali Kamali, Allan Jabri, Allison Moyer, Allison Tam, Amadou Crookes, Amin Tootoochian, Amin Tootoonchian, Ananya Kumar, Andrea Vallone, Andrej Karpathy, Andrew Braunstein, Andrew Cann, Andrew Codispoti, Andrew Galu, Andrew Kondrich, Andrew Tulloch, Andrey Mishchenko, Angela Baek, Angela Jiang, Antoine Pelisse, Antonia Woodford, Anuj Gosalia, Arka Dhar, Ashley Pantuliano, Avi Nayak, Avital Oliver, Barret Zoph, Behrooz Ghorbani, Ben Leimberger, Ben Rossen, Ben Sokolowsky, Ben Wang, Benjamin Zweig, Beth Hoover, Blake Samic, Bob McGrew, Bobby Spero, Bogo Giertler, Bowen Cheng, Brad Lightcap, Brandon Walkin, Brendan Quinn, Brian Guaraci, Brian Hsu, Bright Kellogg, Brydon Eastman, Camillo Lugaressi, Carroll Wainwright, Cary Bassin, Cary Hudson, Casey Chu, Chad Nelson, Chak Li, Chan Jun Shern, Channing Conger, Charlotte Brette, Chelsea Voss, Chen Ding, Cheng Lu, Chong Zhang, Chris Beaumont, Chris Hallacy, Chris Koch, Christian Gibson, Christina Kim, Christine Choi, Christine McLeavey, Christopher Hesse, Claudia Fischer, Clemens Winter, Coley Czarnecki, Colin Jarvis, Colin Wei, Constantin Koumouzelis, and Dane Sherburn et al. (318 additional authors not shown). Gpt-4o system card, 2024. URL <https://arxiv.org/abs/2410.21276>.
- [17] OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutscher, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, and Shawn Jain et al. (181 additional authors not shown). Gpt-4 technical report, 2024. URL <https://arxiv.org/abs/2303.08774>.

- [18] Yujia Qin, Yining Ye, Junjie Fang, Haoming Wang, Shihao Liang, Shizuo Tian, Junda Zhang, Jiahao Li, Yunxin Li, Shijue Huang, et al. Ui-tars: Pioneering automated gui interaction with native agents. *arXiv preprint arXiv:2501.12326*, 2025.
- [19] Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy Lillicrap. Android in the wild: A large-scale dataset for android device control, 2023. URL <https://arxiv.org/abs/2307.10088>.
- [20] John Steven, Pravir Chandra, Bob Fleck, and Andy Podgurski. jrapture: A capture/replay tool for observation-based testing. In *Proceedings of the 2000 ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA ’00*, page 158–167, New York, NY, USA, 2000. Association for Computing Machinery. ISBN 1581132662. doi: 10.1145/347324.348993. URL <https://doi.org/10.1145/347324.348993>.
- [21] Liangtai Sun, Xingyu Chen, Lu Chen, Tianle Dai, Zichen Zhu, and Kai Yu. Meta-gui: Towards multi-modal conversational agents on mobile gui. *arXiv preprint arXiv:2205.11029*, 2022.
- [22] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy Lillicrap, Angeliki Lazaridou, Orhan Firat, James Molloy, Michael Isard, Paul R. Barham, Tom Hennigan, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, Ryan Doherty, Eli Collins, Clemens Meyer, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel, Jack Krawczyk, Cosmo Du, Ed Chi, Heng-Tze Cheng, Eric Ni, Purvi Shah, Patrick Kane, Betty Chan, Manaal Faruqui, Aliaksei Severyn, Hanzhao Lin, YaGuang Li, Yong Cheng, Abe Ittycheriah, Mahdis Mahdieh, Mia Chen, Pei Sun, Dustin Tran, Sumit Bagri, Balaji Lakshminarayanan, Jeremiah Liu, Andras Orban, Fabian Güra, Hao Zhou, Xinying Song, Aurelien Boffy, Harish Ganapathy, Steven Zheng, HyunJeong Choe, Ágoston Weisz, Tao Zhu, Yifeng Lu, Siddharth Gopal, Jarrod Kahn, Maciej Kula, Jeff Pitman, Rushin Shah, Emanuel Taropa, Majd Al Merey, Martin Baeuml, Zhifeng Chen, Laurent El Shafey, Yujing Zhang, Olcan Sercinoglu, George Tucker, Enrique Piqueras, Maxim Krikun, Iain Barr, Nikolay Savinov, Ivo Danihelka, Becca Roelofs, Anaïs White, Anders Andreassen, Tamara von Glehn, Lakshman Yagati, Mehran Kazemi, Lucas Gonzalez, Misha Khalman, Jakub Sygnowski, Alexandre Frechette, Charlotte Smith, Laura Culp, Lev Proleev, Yi Luan, and Xi Chen et al. (1250 additional authors not shown). Gemini: A family of highly capable multimodal models, 2024. URL <https://arxiv.org/abs/2312.11805>.
- [23] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.
- [24] Jialong Wu, Wenbiao Yin, Yong Jiang, Zhenglin Wang, Zekun Xi, Runnan Fang, Deyu Zhou, Pengjun Xie, and Fei Huang. Webwalker: Benchmarking llms in web traversal, 2025. URL <https://arxiv.org/abs/2501.07572>.
- [25] 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, 2024.
- [26] Henry Hengyuan Zhao, Difei Gao, and Mike Zheng Shou. Worldgui: Dynamic testing for comprehensive desktop gui automation, 2025. URL <https://arxiv.org/abs/2502.08047>.

## A Dataset Format

Each datapoint is stored as a single JSON file accompanied by a set of PNG screenshots, one for each interaction step. The JSON file contains the following fields:

- **task\_description:** A natural language instruction describing the overall task goal the agent should complete.
- **A set of Steps:** Each task consists of multiple interaction steps, each containing the following fields:
  - **step\_description:** A natural language description of the expected action at this step, grounded in the corresponding screenshot.
  - **img\_path:** The file name to the screenshot image (PNG format) representing the current GUI state.
  - **ele\_loc/x:** The x-coordinate of the top-left corner of the target UI element.
  - **ele\_loc/y:** The y-coordinate of the top-left corner of the target UI element.
  - **ele\_loc/width:** The width of the bounding box surrounding the target UI element.
  - **ele\_loc/height:** The height of the bounding box surrounding the target UI element.
  - **ele\_type:** The type of the UI element, categorized as icon, text, or box.
  - **action/type:** The type of action to be performed, selected from click, input, get\_info, open, wait, or human.
  - **action/content:** The semantic content of the action, such as the text to input or the label of a button to click.
  - **action\_human (optional):** For abnormal scenarios, this field provides example content for the human action type, describing the issue that the agent should report to a human operator.

Table 5 defines the complete action space used in GUI-Robust. Each step in a task is annotated with an action composed of two components: a type, which specifies the kind of user interaction (e.g., click, input), and a content, which provides semantic detail about the interaction (e.g., the text to input or the target element to click):

- **click:** Represents a standard mouse click operation. The content field describes the semantic label or textual content of the clicked element (if available), such as a clickable text labeled “Search” or “Next”.
- **input:** Denotes a text input operation into a form field or search bar. The content provides the actual text string to be entered, such as “comic book” in a search box.
- **get\_info:** Indicates an information retrieval action. The content describes the specific piece of information to be extracted from the GUI, such as a price, a user name, or a product description.
- **wait:** A passive action used when the GUI is temporarily unresponsive (e.g., during page loading). This type has no content field, as the action itself signals a need to pause before proceeding.
- **open:** Represents an explicit instruction to launch a new web page or desktop application. The content field specifies the name or identifier of the target app or site to be opened.
- **human:** Used when the agent encounters an abnormal scenario it cannot resolve autonomously, such as login screens or Captcha challenges. The content field provides a human-readable description of the issue, serving as a report for fallback intervention.

Table 5: Action space(type + content) and their semantic meanings

Type	Content	Description
click	The content of the clicked element (if any)	A click action.
input	The text to be entered.	A text input action.
get_info	The information be retrieved.	An operation to retrieve information from the interface;
wait	None	Signifies a wait action
human	A description of the problem encountered by agent.	An action to report the abnormal scenario.
open	The name of the web/app to open.	Action to open a web/app

## B Data Example

See an example in Fig. 3

**Task Description:** Search for the best-selling comic book on Amazon

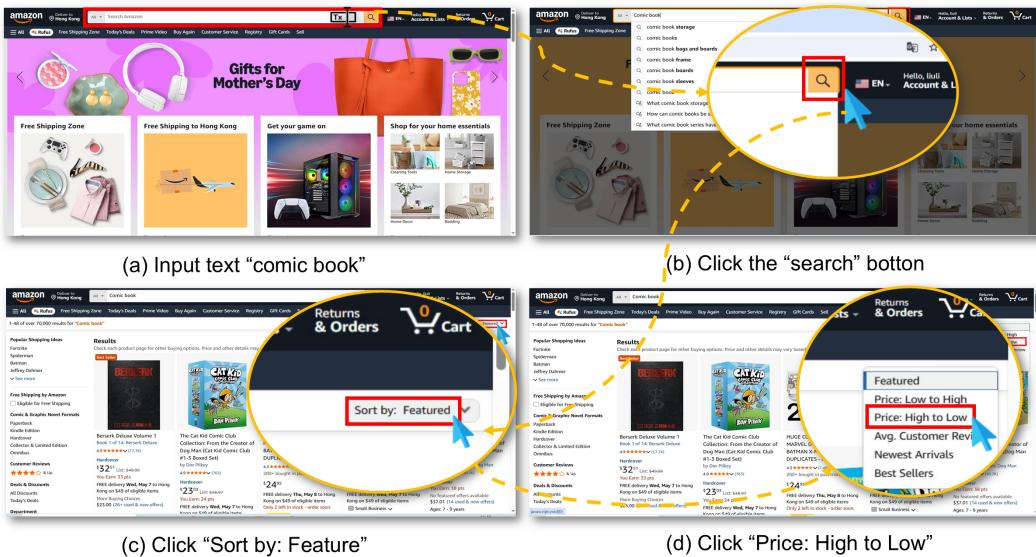


Figure 3: An example task from GUI-Robust: Search for the best-selling comic book on Amazon. The task is decomposed into four steps: (a) enter the query “comic book” in the search bar, (b) click the search button, (c) open the sorting dropdown, and (d) select the “Price: High to Low” option. Each step corresponds to a user interaction grounded in the GUI screenshot.

Table 6: Abnormal Scenarios and Expected Actions

Abnormal Scenario	Expected Action
Action Failure	1. human: Detected failure of the previous action; 2. click/input (repeat previous action)
Login Page	human: Detected the login page, which needs to be completed manually
Captcha Page	human: Detected the captcha page, which needs to be completed manually
Ad Pop-up	1. human: Detected an ad pop-up blocking the process; 2. click (close the pop-up)
Cookie Pop-up	1. human: Detected a cookie pop-up blocking the process; 2. click (accept cookies)
Page Loading	1. human: Detected that the page is loading; 2. wait
Network Disconnection	human: Detected that the network is currently disconnected

## C Abnormal Scenarios Example

In this appendix, we provide illustrative examples and formal definitions of abnormal GUI scenarios included in the GUI-Robust dataset. These scenarios reflect common disruptions that frequently occur in real-world GUI usage but are largely overlooked in prior benchmarks.

Fig. 4 presents visual examples of nine representative abnormal conditions. These include transient interface elements such as cookie consent dialogs (a) and advertisement pop-ups (b), as well as system-level disruptions such as login prompts (c, d), Captcha challenges (g–i), page loading delays (e), and network disconnections (f). Each of these anomalies can significantly interrupt the agent’s execution flow and requires appropriate detection and response strategies.

Table 6 summarizes the expected agent behaviors in response to each abnormal scenario. For instance, when encountering a login or Captcha page, the agent should generate a human action to report that manual intervention is required. In contrast, if a pop-up can be resolved autonomously (e.g., cookie acceptance or ad dismissal), the agent should attempt to issue a valid `click` action on the relevant UI element. In some cases, such as delayed page loading, a `wait` action is necessary before retrying the intended operation.

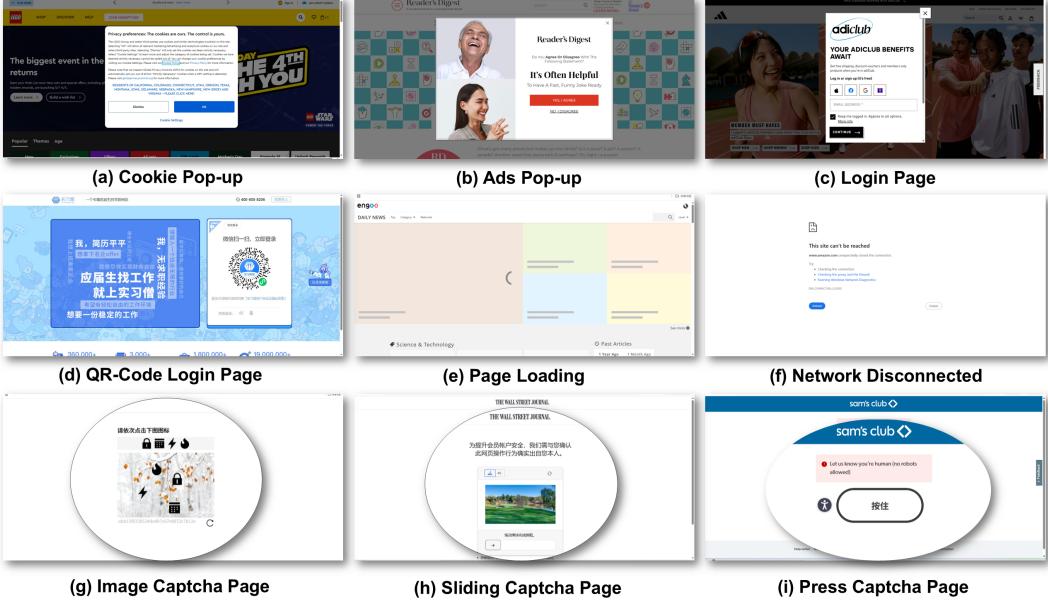


Figure 4: Examples of abnormal GUI scenarios in the GUI-Robust dataset. (a) cookie pop-up, (b) advertisement pop-up, (c) login page, (d) QR-code login prompt, (e) page loading state, (f) network disconnection, (g) image-based Captcha, (h) sliding Captcha, and (i) press-to-verify Captcha. These scenarios require agents to detect and appropriately handle interruptions in order to maintain task robustness.

## D Toolkit Integration Guide

This appendix provides details on how to integrate custom models into the evaluation pipeline.

**Model Interface Specification** Each model should implement the following interface:

```
class YourModel:
    def __init__(self):
        pass

    def pred_step_loc(self, step_description: str, screenshot_base64: str) \
            -> dict:
        # Return a prediction for the current step
        return {
            "ele_loc": {"x": 100, "y": 200},
            "ele_type": "text",
            "action": {"type": "click", "content": "Search"}
        }

    def pred_task_full(self, task_description: str, \
                      screenshot_list_base64: List[str]) -> List[dict]:
        # Return a list of predictions for the full task
        return [ ... ] # One entry per step
```

These two methods respectively handle single-step prediction and full-task multi-step prediction. The evaluation pipeline will automatically invoke these methods, compare predictions against ground truth, and report metrics such as action accuracy, coordinate accuracy, and task success rate.

**Prediction Output Format** Each prediction should contain:

- Element coordinates (x, y)
- Element type (icon, text, box, none)

- Action type and content (click, input, get\_info, open, wait, human)

For full-task predictions, return a list of such dictionaries (one per step).

**Evaluation Script** Run the toolkit via:

```
python evaluation.py \
--model_name YourModel \
--eval_type step|task \
--task_type normal|abnormal \
--data_path path_to_data_folder
```

The evaluation script supports two modes:

-eval\_type step: Evaluates per-step grounding accuracy (Action acc. and Coord. acc.)

-eval\_type task: Evaluates full task execution (Action acc., Coord. acc., and Task Success)

**Prompt Templates** The evaluation toolkit internally uses task-specific prompt templates for multi-modal models. For example:

```
You are a GUI agent. You are given a step with screenshot.
You need to perform the next action to complete the task.

## Output Format
```
{
  "ele_loc": "(x, y)",
  "ele_type": "",
  "action": {
    "type": "",
    "content": ""
  }
}
```

## action Space

| type | content | Description |
| ----- | | |
| click | The content of the element clicked (if any), for example, when clicking on a text element, content is the text content | |
| click | click |
| input | Input text | Input text |
| get_info | Obtained information | Information acquisition operation, some operations involve obtaining information in the interface |

## ele_loc Space

| ele_type | Description |
| ----- | |
| icon | Abstract icon class element, that is, an icon without text |
| text | Text class element, that is, an icon containing text |
| box | Elements such as input boxes |
| none | has no element type. The action get_info does not correspond to any UI element, so the type is none |

## Note
- Use Chinese if the 'content' is Chinese. And use English if the 'content' is English.
```

Our codebase is publicly available at:

<https://github.com/chessbean1/GUI-Robust>

## E Details of Manual Data Collection Pipeline

To provide a meaningful baseline for comparison, we also implemented a traditional fully manual data collection pipeline. This five-step approach mirrors common practices used in existing GUI datasets and serves as a contrast to our proposed semi-automated RevAct pipeline. The workflow is as follows:

**Step 1: Task Design.** Annotators first select a specific application or website and design a meaningful task based on real-world usage scenarios (e.g., "check the latest train schedule", "send a message via the chat app").

**Step 2: Step Decomposition.** Each task is decomposed into atomic steps, corresponding to the minimal sequence of user interactions required to complete the task. Each step should reflect a single action on a clearly identifiable UI element.

**Step 3: Screenshot Capture.** For every step, annotators manually take a screenshot of the GUI interface immediately before the interaction occurs. These screenshots serve as the visual grounding context for the agent.

**Step 4: Annotation.** Annotators label each step with the following fields: a natural language step description, element location (bounding box coordinates), element type, and action type/content. For abnormal tasks, the corresponding human action content is also specified.

**Step 5: Review and Correction.** A second annotator independently checks each task for correctness and consistency. Common issues include incorrect coordinates, ambiguous step descriptions, or wrong actions. Such issues are manually flagged and revised to ensure annotation quality.

In practice, we observed several challenges associated with the traditional manual data collection process, which significantly limit its scalability and consistency. These include:

- **High Learning Overhead for New Applications:** Annotators often face steep learning curves when unfamiliar with the target software. Understanding the interface structure, navigation flow, and interaction logic may require extensive exploration, increasing the overall collection time and cognitive burden.
- **Annotation Complexity and Quality Sensitivity:** In addition to executing tasks, annotators must carefully label each step with precise formatting. They must also pay attention to dynamic interface states, such as pop-up dialogs, cookie banners, or login modals, which may partially or fully obscure the target UI elements.
- **Human Variability and Subjective Interpretation:** Different annotators may interpret the same task or GUI state differently, leading to inconsistencies in task decomposition, step description, or element labeling. Without strict guidelines and training, this can introduce annotation noise and reduce benchmark reliability.
- **Error Propagation and Correction Overhead:** Errors made early in the process—such as a misinterpreted task goal or incorrect screenshot timing—can propagate through all subsequent steps. Detecting and fixing these mistakes requires full reprocessing of the trajectory, making error correction costly and inefficient.

These limitations highlight the practical difficulties of large-scale, high-quality GUI annotation when relying solely on manual pipelines. They further motivate the need for semi-automated approaches such as RevAct that can reduce human burden while preserving annotation richness and realism.

## NeurIPS Paper Checklist

### 1. Claims

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

Answer: [Yes]

Justification: The contribution of this paper is summarized in the end of Introduction.

Guidelines:

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

### 2. Limitations

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

Answer: [Yes]

Justification: We discuss the limitation of this paper in Section 6

Guidelines:

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

### 3. Theory assumptions and proofs

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

Answer: [NA]

Justification: Our work focuses on a new dataset construction, therefore does not include theoretical results.

Guidelines:

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

#### 4. Experimental result reproducibility

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

Answer: [Yes]

Justification: We have released the evaluation scripts on GitHub, and provide detailed usage instructions in Appendix D.

Guidelines:

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

#### 5. Open access to data and code

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

Answer: [Yes]

Justification: We have uploaded the code at <https://github.com/chessbean1/GUI-Robust>

Guidelines:

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

## 6. Experimental setting/details

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

Answer: [Yes]

Justification: The detailed experimental settings are introduced in section 5.1 and Appendix D.

Guidelines:

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

## 7. Experiment statistical significance

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

Answer: [No]

Justification: We report deterministic performance metrics for all models based on fixed test sets without repeated trials or variance analysis; thus, no statistical significance or error bars are provided.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments compute resources

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

Answer: [\[Yes\]](#)

Justification: All experiments were conducted on A100 GPUs (40GB), using either 1 or 2 GPUs depending on model size (e.g., UI-TARS-7B vs. Qwen2.5-VL-72B). Each evaluation run on our benchmark typically took 20–60 minutes per model. We estimate the total compute budget for all reported experiments to be approximately 150 GPU-hours.

Guidelines:

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

## 9. Code of ethics

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

Answer: [\[Yes\]](#)

Justification: Our research complies with the NeurIPS Code of Ethics. All data were collected from publicly accessible platforms with no personally identifiable information.

Guidelines:

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

## 10. Broader impacts

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

Answer: [\[Yes\]](#)

Justification: Our dataset supports the development of more robust and generalizable GUI agents, which could reduce human workload in software operation, accessibility, and automation.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

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

Answer: [NA]

Justification: Our work does not involve releasing models or data with high risk of misuse. The dataset contains only GUI interaction traces from publicly accessible applications and does not include sensitive content or private user information.

Guidelines:

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

## 12. Licenses for existing assets

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

Answer: [Yes]

Justification: All external assets used in this work—such as pretrained models (e.g., Qwen2.5-VL) and benchmarks (e.g., Mind2Web, WebLINX)—are properly cited with references to their original papers. We adhere to their respective open-source licenses (e.g., Apache 2.0, CC-BY 4.0) where applicable, and no proprietary or restricted resources were used.

Guidelines:

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

### 13. New assets

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

Answer: [Yes]

Justification: We release a new dataset (GUI-Robust) and evaluation scripts, both of which are accompanied by detailed documentation.

Guidelines:

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

### 14. Crowdsourcing and research with human subjects

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

Answer: [No]

Justification: Our data collection involved in-house annotators rather than public crowdsourcing platforms. While we provided task guidelines and interface demonstrations, we did not include the full instruction text or compensation details in the paper or supplementary material.

Guidelines:

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

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

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

Answer: [Yes]

Justification: Our study does not involve research with human subjects as defined by IRB standards. All data were collected via non-identifiable interaction traces from internal annotators performing routine tasks on public applications.

Guidelines:

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

## 16. Declaration of LLM usage

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

Answer: [Yes]

Justification: Our data collection pipeline (RevAct) leverages the Qwen2.5-VL multimodal large language model to generate step-wise and task-level natural language descriptions based on user interaction traces. This usage is integral to the core methodology and is described in detail in Section 4.

Guidelines:

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