

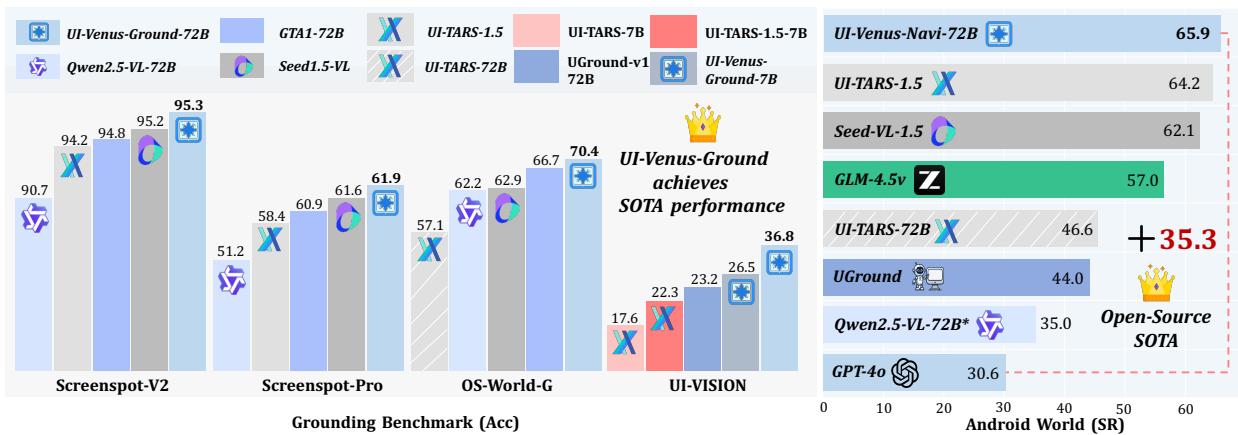
# UI-Venus Technical Report: Building High-performance UI Agents with RFT

Zhangxuan Gu\*, Zhengwen Zeng\*, Zhenyu Xu\*, Xingran Zhou\*, Shuheng Shen\*,†, Yunfei Liu\*, Beiting Zhou\*, Changhua Meng, Tianyu Xia, Weizhi Chen, Yue Wen, Jingya Dou, Fei Tang, Jinzhen Lin, Yulin Liu, Zhenlin Guo, Yichen Gong, Heng Jia, Changlong Gao, Yuan Guo, Yong Deng, Zhenyu Guo, Liang Chen, Weiqiang Wang

Ant Group

We present UI-Venus, a native UI agent that takes only screenshots as input based on a multimodal large language model. UI-Venus achieves SOTA performance on both UI grounding and navigation tasks using only several hundred thousand high-quality training samples through reinforcement finetune (RFT) based on Qwen2.5-VL. Specifically, the 7B and 72B variants of UI-Venus obtain **94.1% / 50.8%** and **95.3% / 61.9%** on the standard grounding benchmarks, *i.e.*, Screenspot-V2 / Pro, surpassing the previous SOTA baselines including open-source GTA1 and closed-source UI-TARS-1.5. To show UI-Venus’s summary and planing ability, we also evaluate it on the AndroidWorld, an online UI navigation arena, on which our 7B and 72B variants achieve **49.1%** and **65.9%** success rate, also beating existing models. To achieve this, we introduce carefully designed reward functions for both UI grounding and navigation tasks and corresponding efficient data cleaning strategies. To further boost navigation performance, we propose Self-Evolving Trajectory History Alignment & Sparse Action Enhancement that refine historical reasoning traces and balances the distribution of sparse but critical actions, leading to more coherent planning and better generalization in complex UI tasks. Our contributions include the publish of SOTA open-source UI agents, comprehensive data cleaning protocols and a novel self-evolving framework for improving navigation performance, which encourage further research and development in the community.

Code: <https://github.com/antgroup/UI-Venus>



**Figure 1** UI-Venus achieves SOTA performance across multiple UI grounding and navigation benchmarks.

\*Equal contribution. †Corresponding author: Shuheng Shen(shuheng.ssh@antgroup.com).

# 1 Introduction

Recent studies of multimodal large language models (MLLMs) Bai et al. (2023); Anthropic (2024); Wang et al. (2024d); Bai et al. (2025); Zhu et al. (2025); Zhipu-AI (2025) have contributed significantly to the advancements and developments of UI agents Hu et al. (2024); Gao et al. (2024); Wang et al. (2024e); Nguyen et al. (2024); Zhang et al. (2024a). In particular, many early approaches, *e.g.*, CogAgent Hong et al. (2024) and UI-TARS Qin et al. (2025), directly leverage extensive open-source and private datasets through pretraining and supervised fine-tuning (SFT) to achieve commendable performance in UI agent tasks. Specifically, these pretrained and SFT approaches treat the complete UI task traces, state observations, and current actions as plain texts for token-level supervised learning. Although SFT is effective in many generation tasks with strong instruction-following ability, sometimes it is not appropriate in discriminative tasks like UI agent. For example, in UI grounding tasks, a predicted point is considered correct as long as it falls within the ground-truth bounding box. In contrast, SFT assigns widely-used cross-entropy loss penalty to any predicted points within and without the box except the center point. Moreover, the data collection and cleaning workflow during pretraining is time consuming and requires a lot of human work.

Inspired by the emergence of DeepSeek-R1 DeepSeek-AI (2025) and its innovative Group Relative Policy Optimization (GRPO) algorithm Shao et al. (2024), recent researchers focus on reinforcement finetune (RFT) paradigms for discriminative tasks such as math and code. RFT generally needs fewer training data, yet has better generalization abilities compared to SFT, as exemplified by Chen et al. (2025). As a result, many approaches similar to UI-R1 Lu et al. (2025b) have also achieved satisfactory results on UI grounding benchmarks by adapting and modifying the reward functions of VLM-R1 Shen et al. (2025). For example, GUI-G1 Zhou et al. (2025), GUI-G<sup>2</sup> Tang et al. (2025a), Phi-Ground Zhang et al. (2025a) and GTA1 Yang et al. (2025) mainly focus on UI grounding benchmarks (ScreenSpotv2 Wu et al. (2024) and ScreenSpotpro Li et al. (2025b)) with different data collection and ratios. Some tricks like think / no-think and input message construction are also discussed in their reports. On the other hand, GUI-R1 Luo et al. (2025a), InfiGUI-R1 Liu et al. (2025) and UI-R1 Lu et al. (2025b) extend GRPO to UI offline navigation benchmarks (Andorid Control Li et al. (2024) and GUI-Odyssey Lu et al. (2025a)) by adapting action type, format and the point rewards with predefined action spaces.

Despite these advancements, current R1-like UI agents exhibit three critical limitations. First, although some attempts are made by UI-R1 Lu et al. (2025b) and InfiGUI-R1 Liu et al. (2025) in offline UI navigation tasks, their summary, memory and planning ability are still not enough for online tasks like AndroidWorld Rawles et al. (2025) when given a user goal with complex and changeful interactive environment. Secondly, one important factor for training UI agents is the data quality. According to our observations, approximately a half of open-source UI data contain noise, while only a few methods have considered the data cleaning and selection strategies. Thirdly, existing implementations mainly focus on small MLLMs (*e.g.*, 3B/7B parameters), neglecting the potential of large-scale models (*e.g.*, 72B) in RFT training. Although GTA1 has developed 72B grounding model, no further attempts are made on the end-to-end UI navigation tasks without addtional planner like GPT4o OpenAI (2024). This constraint results in performance gaps compared to state-of-the-art large-scale models like UI-TARS-1.5 Seed (2025b) and SeedVL-1.5 Seed (2025a) on UI agent benchmarks. Moreover, the evaluation for UI agents on many benchmarks suffers from severe challenges due to different input prompts as well as unreleased hyper-parameters and dataset settings. For example, some results from the official papers can not be reproduced according to Yang et al. (2025); Zhang et al. (2025b) and some github issues.

To address the first challenge mentioned above, we design a Self-Evolving Trajectory History Alignment & Sparse Action Enhancement framework. This method addresses two critical limitations in existing UI navigation agents: (1) misaligned historical reasoning traces and (2) insufficient learning of rare but pivotal actions. For the first issue, we iteratively refine the thought-action histories between training epochs, aligning the form and level of detail of the historical reasoning with the agent’s evolving decision-making patterns. This produces a more coherent and informative context for predicting subsequent steps, which in turn improves planning accuracy. For the second, we selectively re-sample trajectories containing sparse actions, constructing multiple historical context variants that lead to the same low-frequency operation. This balanced sampling increases the model’s exposure to infrequent yet crucial skills, enhancing its ability to generalize in complex and dynamic UI task scenarios.

To acquire high-quality UI data, we implement a three-stage processing pipeline to tackle the second problem: (1) **Data Filtering** includes unifying scroll directions, filtering out trajectories with inconsistencies, and resampling trajectories based on their categories. (2) **Trace Reconstruction** refers to modifying the information-retrieval traces with specific answers inserted. (3) **Iteratively Trace Generation:** we develop a data generation framework by using UI-Venus-Navi to predict and record trajectories, with comprehensive quality filtering strategies to select high-quality traces for iterative training. During the training, we use about 107k/350k high-quality training samples filtered by ourselves for UI grounding/navigation task, respectively. More details will be introduced in the next section.

In this report, we develop and publicly release the UI-Venus series, comprising UI-Venus-Ground-7B/72B and UI-Venus-Navi-7B/72B, all trained with GRPO on the Qwen2.5-VL model [Bai et al. \(2025\)](#). We also open-source the evaluation codes of grounding as well as the prompts and post-processing scripts of navigation to enhance accessibility and ease of use for researchers. Experimental results exhibit that UI-Venus outperforms all existing UI agents, showing strong performance with about 350k self-constructed high quality dataset. Our model achieves new SOTA performance on five UI grounding benchmarks including ScreenSpot-V2 [Wu et al. \(2024\)](#), ScreenSpot-Pro [Li et al. \(2025b\)](#), OSWorld-G [Xie et al. \(2025\)](#), UI-Vision [Nayak et al. \(2025\)](#) and CA-GUI [Zhang et al. \(2025b\)](#). The 7B and 72B variants of UI-Venus obtain **94.1% / 50.8%** and **95.3% / 61.9%** on the Screenspot-V2 / Pro, surpassing previous SOTA baselines including open-source GTA1 and closed-source UI-TARS-1.5. For UI navigation, we evaluate UI-Venus on both offline and online UI navigation benchmarks. Our model achieves SOTA performance on AndroidWorld [Rawles et al. \(2025\)](#), where 7B and 72B variants achieve **49.1%** and **65.9%** success rate, while obtaining comparable results on AndroidControl [Li et al. \(2024\)](#) and GUI-Odyssey [Lu et al. \(2025a\)](#) benchmarks.

There are two motivations for publishing the grounding and navigation models separately. (1) **A More Efficient Grounding Model:** According to our experiments, explicitly prompting the model to output its thinking process significantly enhances its observation and planning capabilities in UI navigation tasks. However, UI-Venus achieves approximate performance on UI grounding tasks with and without thinking. As no-think mode for grounding outputs only several location tokens, it’s more efficient in inference for real-world applications. (2) **Reward Conflicts:** If we attempt to train an agent using both grounding and navigation tasks, a feasible approach is to convert grounding data into single-step click actions in the navigation task. However, using these two reward mechanisms simultaneously can lead to unstable training, ultimately degrading performance on both tasks.

We summarize our main contributions as follows:

- To mitigate historical reasoning misalignment and augment the learning of rare but pivotal actions, we design a Self-Evolving Trajectory History Alignment & Sparse Action Enhancement framework, which in turn boosts navigation performance in complex UI scenarios.
- We conduct a comprehensive study of the quality of UI data, and propose a data cleaning and selection strategy to improve training data quality for both grounding and navigation.
- We develop and open-source UI-Venus, a SOTA UI agent for both grounding and navigation tasks with carefully designed reward functions, whose 7B and 72B variants obtaining **94.1% / 50.8%** and **95.3% / 61.9%** on Screenspot-V2 / Pro, and **49.1%** and **65.9%** on AndroidWorld, demonstrating the effectiveness of model scaling up for GRPO in UI-agent tasks.

## 2 Related Works

### 2.1 UI Grounding

UI Grounding focuses on localizing and identifying UI elements based on natural language instructions, serving as the foundation of automated UI interaction. Traditional approaches have relied on Supervised Fine-Tuning to train models on labeled UI datasets [Cheng et al. \(2024\)](#); [Lin et al. \(2024\)](#); [Xu et al. \(2024\)](#); [Wu et al. \(2024\)](#); [Gu et al. \(2023\)](#); [Gou et al. \(2024\)](#); [Wang et al. \(2024c\)](#). However, these methods face two primary limitations: (1) poor generalization in out-of-distribution scenarios where UI layouts or visual styles differ from training data, and (2) high costs associated with acquiring large-scale annotated datasets for diverse UI environments. Recent advances have shifted toward reinforcement learning-based fine-tuning inspired by the DeepSeek-R1 paradigm. Early works like UI-R1 [Lu et al. \(2025b\)](#) and GUI-R1 [Luo et al. \(2025a\)](#) introduced binary hit-or-miss rewards for task completion, while InfiGUI-R1 [Liu et al. \(2025\)](#) proposed two-stage training combining offline pretraining with online RL. GUI-G1 [Zhou et al. \(2025\)](#) further refined the reward design with box-size-based rewards for improved spatial precision. The latest methods, including SE-GUI [Yuan et al. \(2025\)](#), LPO [Tang et al. \(2025c\)](#), and GUI-G<sup>2</sup> [Tang et al. \(2025a\)](#), employ continuous reward mechanisms that provide fine-grained feedback throughout the grounding process.

### 2.2 UI Agents

**UI Agent Framework.** The UI agent framework leverages collaborative agent systems to handle complex GUI automation tasks through task decomposition and specialization. Mobile-Agent [Wang et al. \(2024b,a, 2025b\)](#) addresses mobile device automation by introducing a planning-decision-reflection architecture, where specialized agents handle task progress tracking, action execution, and operation verification respectively, while maintaining a memory unit for context preservation across interactions. Cradle [Tan et al. \(2024\)](#) presents a modular framework with six core components—information gathering, self-reflection, task inference, skill curation, action planning, and memory—enabling agents to tackle both video game control and software manipulation through adaptive learning and skill reuse. Agent-S [Agashe et al. \(2024\)](#) implements role-specific agents that specialize in distinct UI interaction patterns, improving task execution through modular decision-making. DroidRun [dro \(2025\)](#) focuses on Android automation by leveraging Accessibility Services to access structured UI hierarchies rather than pixel-based approaches, enabling reliable interaction with native mobile applications . These frameworks illustrate that agent specialization, in conjunction with structured communication protocols, can effectively address diverse UI scenarios. However, this capability comes at the cost of increased system complexity and computational overhead required to

coordinate multiple agents.

**Native UI Agent.** Native UI agents Hong et al. (2024); Qin et al. (2025); Zhang et al. (2025b,a); Feng et al. (2025); Bai et al. (2024); Humphreys et al. (2022) represent a paradigm shift toward unified and end-to-end systems that directly learn to interact with graphical interfaces without requiring multiple specialized components. CogAgent Hong et al. (2024) and UI-TARS Qin et al. (2025) pioneered this approach through large-scale training on diverse UI interaction data, enabling the model to develop a comprehensive understanding of GUI patterns across different platforms and applications. By training on millions of UI screenshots paired with action sequences, UI-TARS demonstrated that a single model could effectively handle various tasks. Following this success, the native UI agent introduced architectural improvements that enhance the model’s ability to process high-resolution screenshots while maintaining efficient action prediction, particularly excelling in scenarios with dynamic UI updates and real-time feedback requirements. AgentCPM-GUI Zhang et al. (2025b) took a different approach by focusing on mobile-specific optimizations, incorporating touch gesture understanding and mobile UI design patterns into its training process, resulting in superior performance on mobile platforms with reduced latency. Meanwhile, Phi-Ground Zhang et al. (2025a) advanced the field by integrating sophisticated vision-language alignment techniques, enabling more robust grounding of natural language instructions to visual UI elements through cross-modal attention mechanisms. These native agents benefit from their streamlined architecture and ability to learn complex UI interaction patterns directly from data, though they face challenges in data efficiency and often require substantial computational resources for training on diverse UI environments.

### 3 Methodology

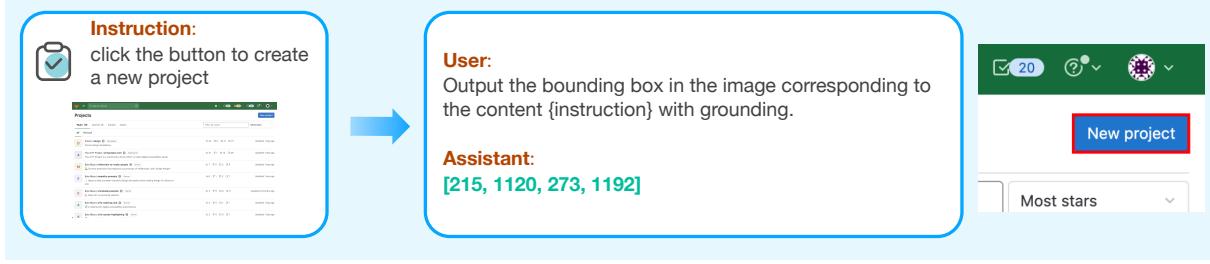
Many prior works have successfully verified that Reinforcement Fine-Tune (RFT) based on Group Relative Policy Optimization (GRPO) algorithm is suitable for UI grounding, where the MLLM answers a response that includes the predicted boxes given corresponding objective descriptions. In this work, we further extend GRPO to UI navigation task and prove its effectiveness in UI Agents training compared to merely SFT. We first present the preliminaries of GRPO, followed by a comprehensive discussion of the data pipelines, reward design, and learning framework in UI grounding and UI navigation respectively.

#### 3.1 Preliminaries

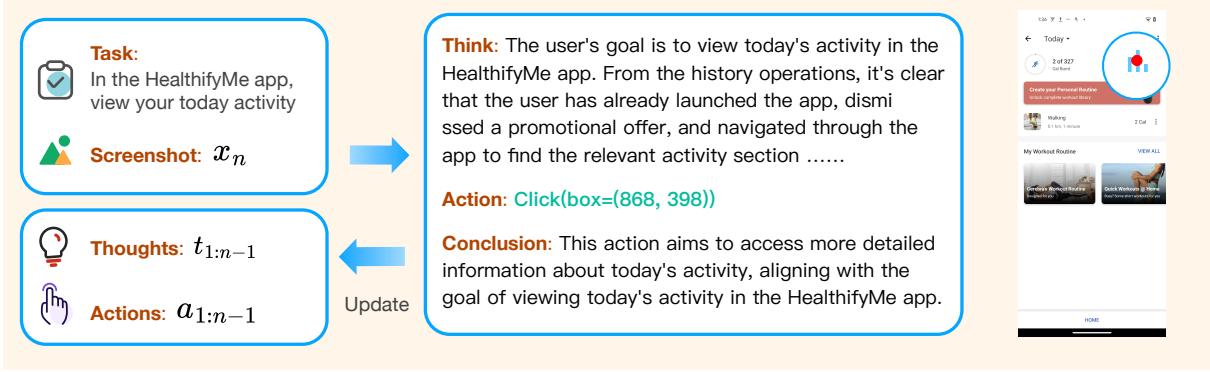
GRPO enhances training stability by estimating baselines through relative rewards within groups rather than using a separate critic model. For each training question  $q \in Q$ , GRPO samples  $G$  rollouts  $\{o_1, o_2, \dots, o_G\}$  from the current policy  $\pi_\theta$  and computes their corresponding rewards  $\{r_1, r_2, \dots, r_G\}$ . The core idea lies in normalizing rewards within each group to obtain advantages:  $\hat{A}_i = \frac{r_i - \text{mean}(\{r_1, r_2, \dots, r_G\})}{\text{std}(\{r_1, r_2, \dots, r_G\})}$ . The policy is then optimized by maximizing the objective:

$$\begin{aligned} \mathcal{J}_{\text{GRPO}}(\pi_\theta) &= \mathbb{E}_{q \sim Q, \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{old}}}(\cdot|q)} \\ &\quad \frac{1}{G} \sum_{i=1}^G \frac{1}{|o_i|} \sum_{t=1}^{|o_i|} \left\{ \min \left[ \frac{\pi_\theta(o_{i,t}|q, o_{i,<t})}{\pi_{\theta_{\text{old}}}(o_{i,t}|q, o_{i,<t})} \hat{A}_i, \text{clip} \left( \frac{\pi_\theta(o_{i,t}|q, o_{i,<t})}{\pi_{\theta_{\text{old}}}(o_{i,t}|q, o_{i,<t})}, 1-\epsilon, 1+\epsilon \right) \hat{A}_i - \beta D_{\text{KL}} [\pi_\theta || \pi_{\text{ref}}] \right] \right\}, \quad (1) \\ &\quad \text{where } \hat{A}_i = \frac{r_i - \text{mean}(\{r_1, r_2, \dots, r_G\})}{\text{std}(\{r_1, r_2, \dots, r_G\})}, \end{aligned}$$

where the clipping mechanism prevents excessive policy updates,  $\epsilon$  controls the clipping range, and the KL divergence term with coefficient  $\beta$  constrains the policy from diverging excessively from the reference model  $\pi_{\text{ref}}$ .



a) UI-Venus-Ground



b) UI-Venus-Navi

**Figure 2** Executions on typical grounding and navigation tasks of UI-Venus. **a)** The instruction and the screenshot are needed for UI-Venus-Ground to output the corresponding coordinates; **b)** For navigation tasks, historical context (thought-action pairs) are essential for UI-Venus-Navi, which will generate the thinking content and model action. The historical context will be updated after each step finished.

## 3.2 UI Grounding

### 3.2.1 Data Collection

Preliminarily, we collect UI grounding instances from existing public datasets, including Widget Captioning [Cheng et al. \(2024\)](#), UI RefExp [Bai et al. \(2021\)](#), SeeClick [Cheng et al. \(2024\)](#), ShowUI [Lin et al. \(2024\)](#), and OmniAct [Kapoor et al. \(2024\)](#). They are extracted from different platforms such as mobile, desktop and web to ensure the diversity of training data. As shown in Table 1, we collect about 627k open-source grounding samples and use about 107k for training after carefully data cleaning.

	Widget Captioning	UI RefExp	SeeClick-Web	ShowUI	OmniAct	Sampled training set
# Samples	34k	17k	325k	110k	141k	<b>107k</b>

**Table 1** GUI grounding data composition. Our training dataset comprises 107k samples from five complementary sources.

### 3.2.2 Data Cleaning

However, according to our observation, approximately 40% of the current open-source grounding data contain significant noise issues, including prompt ambiguity and box shifts. A possible cause is that many web and mobile datasets utilize HTML and A11Y (Accessibility) source code, respectively. Although these source codes contain the bounding boxes and their corresponding descriptive instructions required for grounding tasks, some boxes may have mismatched offsets after page rendering. In severe cases, all boxes in one entire image may shift simultaneously in one certain

direction. Additionally, due to nested elements in UI components from source code, bounding boxes and instructions may fail to maintain a one-to-one correspondence, with many-to-one and one-to-many mappings potentially affecting model training.

Since RFT training requires high-quality clean data, we employ manual inspection to ensure the correctness of box-instruction pairs. Specifically, given the original open source datasets (*e.g.*, Seeclick), which are large in scale but contain many redundant or simple samples, we first perform downsampling to create a subset by removing the repeated prompts. Then we manually filter out ambiguous prompts, relocate offset boxes, and rephrase unmeaningful instructions. After these steps, we obtain about 107k high-quality training samples, which is shown in Table 1.

Failed attempts for automatic grounding data cleaning include: (1) Using open-source models to perform the above filtering and rewriting operations; (2) RFT with hard samples selected by rejection sampling strategies. Although these approaches failed to work, they provide some insight for the UI-agent community.

### 3.2.3 Reward Function

In the UI grounding task, we only use two reward functions, *i.e.*, point-in-box and format reward. Although many approaches Lu et al. (2025b); Tang et al. (2025a); Luo et al. (2025a) use intersection-over-union (IoU) or other smooth rewards, the results are similar according to our experiments. Specifically, the reward function is composed of two components formally:

**Format Reward.** We first check whether the predicted answer string conforms to a predefined syntax. Valid answers receive a base reward to ensure the model to generate executable and parsable instructions.

**Point-in-box Reward.** Given a screenshot and the instruction, the model must predict a bounding box that localizes the element, where  $(x_c, y_c)$  denote the box center. Assume that the ground truth is annotated as  $[x_1, y_1, x_2, y_2]$ , then

$$R_{\text{poin-in-box}} = \begin{cases} 1 & \text{if } x_1 \leq x_c \leq x_2 \text{ and } y_1 \leq y_c \leq y_2, \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

**Total Reward.** Combining all components, the final action-wise reward is computed as:

$$R = R_{\text{format}} \cdot w_1 + R_{\text{poin-in-box}} \cdot w_2, \quad (3)$$

where  $w_1$  and  $w_2$  control the relative importance of format correctness and location precision.

## 3.3 UI Navigation

### 3.3.1 Data Collection

For UI navigation tasks, we first collect traces from open-source navigation datasets of GUI Odyssey Lu et al. (2025a), Aguvis-Aitz Xu et al. (2024), AndroidControl Li et al. (2024), Aitw Rawles et al. (2023a), and Amex Kapoor et al. (2024) as shown in Table 2. Meanwhile, grounding data plays a critical role in improving click operation accuracy for navigation tasks, so we also integrate grounding data into the training process. In addition, to strengthen the model’s ability to handle long-chain traces and capabilities in cross-lingual generalization, we build a custom annotation

Dataset	Multiple Steps						Single Step		Sampled training set
	GUI-O	AITZ	AC	Aitw	Amex	Navi*	GD*	VQA*	
#Samples	114k	14k	84k	371k	39k	20k	107k	42k	<b>350k</b>
#Traces	7k	2k	13k	5k	3k	*	-	-	*

**Table 2** Basic statistics for our collected navigation training data. Note that “AC” means Android Control dataset while the “GUI-O” represents GUI-Odyssey. The datasets marked with \* means are collected by ourselves. Among these data, we select 350k high-quality samples as our training set.

platform that provides approximately 20k samples from several popular Chinese mobile APPs (notated Navi\* in Table 2).

Besides, a unified action space is critical for the integration of data from different sources, enabling the agent to focus on learning actions without being confused by diverse definitions. As shown in Table 3, we follow UI-TARS but modify its action space (*e.g.*, redefine CallUser action for information-retrieval) to better suit the existing open-source navigation training data.

Action	Definition
Click(box=(x, y))	Click at coordinates (x, y).
Drag(start=(x1, y1), end=(x2, y2))	Drag from (x1, y1) to (x2, y2).
Scroll(start=(x1, y1), end=(x2, y2), direction=“)“)	Scroll from (x1, y1) to (x2, y2) with specified direction.
Type(content=“)“)	Type the specified content.
Launch(app=“)“)	Launch the specified app.
Wait()	Wait for loading.
Finished(content=“)“)	Finish the task, with optional information.
CallUser(content=“)“)	Conclude the answer for information-retrieval.
LongPress(box=(x, y))	Long press at coordinates (x, y).
PressBack()	Press the ‘back’ button.
PressHome()	Press the ‘home’ button.
PressEnter()	Press the ‘enter’ button.
PressRecent()	Press the ‘recent’ button.

**Table 3** All actions and their definitions used in UI-Venus. We unify the action space and map all the actions in the existing open-source dataset to this space.

### 3.3.2 Data Pipeline

As noisy data are detrimental to the training process of the model, we develop a pipeline to build a clean training set that can properly guide the model to learn various instructions, which includes three stages: (1) **Data Filtering** to reintegrate existing data, (2) **Trace Reconstruction** to modify existing traces, and (3) **Iteratively Trace Generation** to produce more high-quality traces beyond existing data. After data cleaning, we finally obtain about 350K samples for UI navigation training.

**Data Filtering.** In this stage, we preliminarily filter the collected data by removing overly short traces and standardizing the direction definition of scroll operations among different datasets. In addition, there are some inconsistencies between actions and tasks, *e.g.*, there are traces that may contain more or less operations to complete tasks, or even not follow the requirements of tasks. To filter out these invalid traces, we utilize MLLM to summarize each action and integrate the summaries to obtain an overall description for each trace, which could be compared with the original task. Finally, we categorize the traces based on the apps and subtasks they involve and resample

them to ensure the diversity of the training data, preventing the model from overfitting to specific scenarios without sufficient exploration on diverse tasks.

**Trace Reconstruction.** Among existing datasets, the information retrieval tasks [Rawles et al. \(2025\)](#), an important category of UI navigation tasks, often lack an explicit answer in the final step of the traces. For example, when users make a request of ‘What is the weather today?’ or ‘What is the total price in my shopping cart now?’, they usually expect an answer from the agent, rather than just being navigated to some pages and finding the answer by themselves. To bridge this gap, we select traces of information retrieval tasks from filtered data, and then adopt the MLLM to generate corresponding answers based on the last screenshots of these traces. Finally, we reconstruct the original trace by inserting a CallUser step with generated answer before the Finish step, requiring the agent to report the final answer before ending an information retrieval task.

**Iteratively Trace Generation.** Besides utilizing existing open-source trajectory data, we also design an automated framework to iteratively generate high-quality traces built upon the virtual cloud environment, which contains dozens of available mobiles. We adopt our well-trained UI-Venus model to generate trajectories on real-world tasks including hand-designed and MLLM-generated instructions to ensure their diversity in Chinese and English mobile applications. The noisy and invalid traces would be discarded by combining following steps: (1) **Rule-based filtering**: empirical rules are defined to remove typical errors (*e.g.*, short trajectory with abnormal exit, the repeating invalid actions); (2) **ORM-based filtering**: we trained an outcome reward model (ORM) regarding the whole trajectory as a single example to score on the generated traces and remove those with low scores; (3) **Annotator-based filtering**: Annotators are required to strictly select correct traces from the remaining ones. Also, to learn from failures, the fault actions would be fixed by annotators and added to the train set with its valid trace prefix. The above process helps iteratively optimize our model for complex real-world scenarios.

### 3.3.3 Trajectory History Alignment and Sparse Action Enhancement

Accurate historical context is critical for successful task execution in dynamic UI navigation planning. By leveraging historical information, the agent model can understand the steps already taken, perceive the state transitions [Qin et al. \(2025\)](#); [Chae et al. \(2025\)](#); [Sun et al. \(2025c\)](#), and engage in self-reflection [Shinn et al. \(2023\)](#); [Wu et al. \(2025a\)](#); [Li et al. \(2025c\)](#).

In our approach, the historical context is provided as thought-action pairs, where the thought represents the reasoning process behind taking a particular action at each step, and the action is the actual executed operation. Specifically, when predicting the  $n$ -th step, historical context from the previous  $n - 1$  steps is formulated as:

$$H_{n-1} = [(t_1, a_1), (t_2, a_2), \dots, (t_{n-1}, a_{n-1})]. \quad (4)$$

The context  $H_{n-1}$ , together with the task description and the current UI screenshot, forms the input to the agent model for the  $n$ -th step prediction.

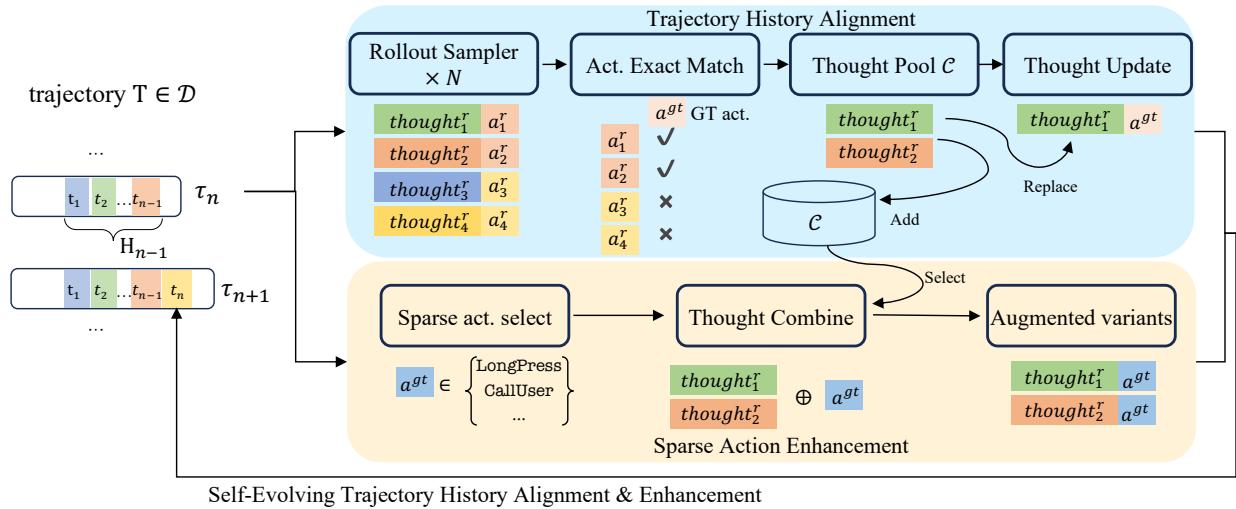
However, in practice, the historical thoughts  $(t_1, \dots, t_{n-1})$  are often not well aligned with the agent model’s intrinsic reasoning capability [Li et al. \(2025a\)](#), due to inconsistencies in style, detail, and abstraction level across data. This mismatch reduces the utility of  $H_{n-1}$ , sometimes causing confusion in decision-making.

To address this issue, we propose a self-evolving [Tao et al. \(2024\)](#) history alignment mechanism that

refines historical thoughts between training epochs, progressively aligning the reasoning traces (*i.e.*, the thought sequence in  $H_{n-1}$ ) with the model’s evolving decision patterns. By dynamically adjusting the trajectory history to better reflect the model’s current reasoning behavior, the mechanism provides a more coherent and consistent historical context, resulting in sustained improvements in navigation performance.

Additionally, we observe that the action distribution in the training data significantly impacts the model performance Qi et al. (2025). Particularly for sparse actions (*e.g.*, LongPress), the agent model often struggles to learn these actions effectively due to their low frequency in the data. In certain complex trajectories, these actions can be pivotal, and missing them may cause the entire task to fail. To overcome the limitation, our method also enhances the sampling of sparse actions, allowing the agent to better acquire and generalize these rare but critical skills. The sampling strategy helps the model develop a more transferable understanding of both the UI states and the planned actions, which improves its robustness in complex and dynamic task scenarios.

To elaborate on these enhancements, we now provide a detailed description of the two key components of our approach. Figure 3 illustrates the overall framework of the proposed *Self-Evolving Trajectory History Alignment & Enhancement* method.



**Figure 3** The overview of the proposed Self-Evolving Trajectory History Alignment & Enhancement process, applied between training epochs. The process consists of two key components: 1) **Trajectory History Alignment** refines the historical context for each trajectory step. The model executes multiple rollouts to generate candidate thought–action pairs, applying an Action Exact Match filter to retain only those whose predicted actions match ground-truth actions. The corresponding thoughts are collected in pool  $\mathcal{C}$  and subsequently replace the original thoughts in the historical context, creating an optimized trajectory history for the next training epoch. 2) **Sparse Action Enhancement** focuses on samples with sparse actions. Multiple variants are constructed by combining different rollout generated thoughts that lead to the same sparse action, effectively increasing the representation of these sparse but critical operations in the training distribution.

**Trajectory History Alignment.** In our framework, the agent model leverages historical context composed of thought–action pairs from previous steps, where each thought describes the reasoning process behind the corresponding action. Explicitly incorporating the past chain-of-thought (CoT) Wei et al. (2023) in this historical context can improve planning quality in sequential decision-making tasks by revealing the latent reasoning behind each action. However, existing UI navigation datasets with high-quality CoT labels are both scarce and heterogeneous. These annotations are collected from diverse sources, including crowd-sourced workers Rawles et al. (2023b), expert

demonstrations [Zhang et al. \(2024b\)](#), and synthetic generation [Sun et al. \(2025b\)](#), which leads to substantial variation in sequence length, linguistic style, and level of reasoning detail.

As historical thoughts play a crucial role in guiding the agent model’s planning, cross-source variability can introduce inconsistencies in its perception for the navigation history, ultimately degrading decision quality. In particular, shorter thought sequences may miss critical intermediate observations, while overly verbose ones may obscure key decision-making milestones.

To address these challenges, we introduce a self-evolving trajectory history alignment mechanism that iteratively refines historical thoughts during training. The detailed procedure is presented in Algorithm 1.

After each training epoch, we perform global trajectory refinement by re-inferring reasoning traces with the current model. Let  $\mathcal{D} = \{T_k\}_{k=1}^K$  denote the dataset. Each trajectory  $T_k$  is a sequence of steps  $\{(x_{k,n}, H_{k,n-1})\}_{n=1}^{N_k}$ , where  $x_{k,n}$  is the UI state (screenshot) at step  $n$ ,  $H_{k,n-1} = [(t_{k,1}, a_{k,1}), \dots, (t_{k,n-1}, a_{k,n-1})]$  is the historical context from the previous  $n - 1$  steps.

For each step  $n$ , given  $(x_{k,n}, H_{k,n-1})$ , we perform  $R$  rollouts:

$$\{(t_{k,n}^{(r)}, a_{k,n}^{(r)})\}_{r=1}^R \sim p_\theta(t, a | x_{k,n}, H_{k,n-1}). \quad (5)$$

We then filter the results to keep only those whose predicted action matches the ground-truth action:

$$\mathcal{C}_{k,n} = \{t_{k,n}^{(r)} \mid a_{k,n}^{(r)} = a_{k,n}^{(gt)}\}. \quad (6)$$

All candidates in  $\mathcal{C}_{k,n}$  form the thought pool  $\mathcal{C}$ .

For the  $n$ -th step, the thoughts in all previous  $n - 1$  steps are refined by selecting a replacement from the corresponding thought pool:

$$t'_{k,i} = \begin{cases} \text{Select}(\mathcal{C}_{k,i}), & \mathcal{C}_{k,i} \neq \emptyset, \\ t_{k,i}, & \text{otherwise,} \end{cases} \quad i = 1, \dots, n-1, \quad (7)$$

where  $\text{Select}(\cdot)$  denotes a selection policy (*e.g.*, choosing the candidate with length closest to a target length).

The updated historical context  $H'_{k,n-1} = [(t'_{k,1}, a_{k,1}^{(gt)}), \dots, (t'_{k,n-1}, a_{k,n-1}^{(gt)})]$  is then used as input in the next training epoch. This refinement is applied after every epoch, aligning the reasoning traces with the agent’s evolving policy. This self-evolving process enhances the agent model’s planning, leading to more robust and coherent navigation performance.

**Sparse Action Enhancement.** The distribution of actions in the training data is often imbalanced (see Figure 4). Common actions such as Click and Scroll appear frequently, while rare actions such as LongPress are much less represented. This imbalance makes it difficult for the model to learn sparse actions effectively. However, these sparse actions often play a critical role in complex task completion, and insufficient handling of such actions results in incomplete or erroneous navigation plans.

To address this problem, we design a sampling strategy that increases learning for sparse actions. During trajectory history alignment, the agent produces a thought pool  $\mathcal{C}_n$  at each step via rollouts. For a step  $(x_n, H_{n-1})$  involving a sparse action (*e.g.*,  $a_n \in \{\text{LongPress}, \text{CallUser}, \dots\}$ ), we create  $M$  historical variants by combining thoughts from the pools. Specifically, for each previous step  $i = 1, \dots, n-1$ , we sample  $\{t'_i^{(m)}\}_{m=1}^M \subset \mathcal{C}_i$  and construct  $H_{n-1}^{(m)} = [(t'_1^{(m)}, a_1), \dots, (t'_{n-1}^{(m)}, a_{n-1})]$ .

---

**Algorithm 1** Self-Evolving Trajectory History Alignment

---

**Require:** Dataset  $\mathcal{D} = \{T_k\}_{k=1}^K$ ; each trajectory  $T_k$  is a sequence of steps  $\{(x_{k,n}, H_{k,n-1})\}_{n=1}^{N_k}$ ; ground-truth actions  $\{a_{k,1:N_k}\}$ ; model  $p_\theta(t, a | x, H)$ ; rollout count  $R$

**Ensure:** Updated historical thoughts  $\{t_{k,1:N_k}\}$

```

1: for each training epoch do
2:   for each trajectory  $T_k$  do
3:     Initialize thought pools  $\{\mathcal{C}_{k,i} \leftarrow \emptyset\}_{i=1}^{N_k}$ 
4:     for  $n = 1$  to  $N_k$  do                                 $\triangleright$  collect candidates for every step
5:        $H_{k,n-1} \leftarrow [(t_{k,1}, a_{k,1}), \dots, (t_{k,n-1}, a_{k,n-1})]$ 
6:       for  $r = 1$  to  $R$  do                       $\triangleright$  rollouts for step  $n$ 
7:          $(t_{k,n}^{(r)}, a_{k,n}^{(r)}) \sim p_\theta(t, a | x_{k,n}, H_{k,n-1})$ 
8:         if  $a_{k,n}^{(r)} = a_{k,n}^{(gt)}$  then           $\triangleright$  action exact-match
9:            $\mathcal{C}_{k,n} \leftarrow \mathcal{C}_{k,n} \cup \{t_{k,n}^{(r)}\}$        $\triangleright$  add to thought pool
10:        end if
11:      end for
12:    end for
13:    for  $i = 1$  to  $N_k$  do                   $\triangleright$  select historical thoughts after collecting pools
14:      if  $\mathcal{C}_{k,i} \neq \emptyset$  then
15:         $t'_{k,i} \leftarrow \text{SELECT}(\mathcal{C}_{k,i})$            $\triangleright$  select policy
16:      end if
17:    end for
18:    for  $n = 1$  to  $N_k$  do                 $\triangleright$  update historical context for next epoch
19:       $H_{k,n-1} \leftarrow [(t'_{k,1}, a_{k,1}^{(gt)}), \dots, (t'_{k,n-1}, a_{k,n-1}^{(gt)})]$ 
20:    end for
21:  end for
22: end for

```

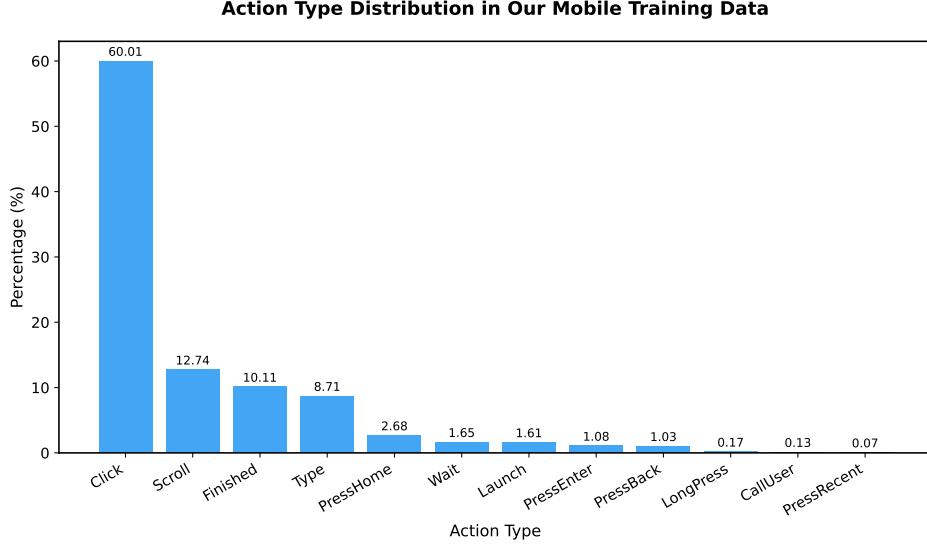
---

Conceptually, the augmented set corresponds to the Cartesian product  $\mathcal{C}_1 \times \dots \times \mathcal{C}_{n-1}$ . This increases both the diversity and frequency of reasoning patterns that lead to sparse actions. For common actions, we do not perform this augmentation. By focusing only on sparse actions, we enhance the model’s ability to learn these rare behaviors, leading to more accurate navigation reasoning.

### 3.3.4 Action-wise Reward Function

The rule-based reward function is a key component in reinforcement learning, especially for tasks where outputs can be verified against well-defined rules. In visual grounding, for example, the reward is often based on the intersection-over-union between the predicted and the ground-truth bounding boxes [Shen et al. \(2025\)](#). While effective for coarse-level correctness, such rewards provide limited guidance for fine-grained UI manipulation.

In UI navigation tasks, the agent model is required to select the correct action type (*e.g.*, Click, Scroll) and generate precise action parameters (*e.g.*, target coordinates or input text). This requirement becomes particularly important in complex interfaces, where even small parameter errors can cause task failure. We propose an *action-wise reward function* tailored for the navigation task, which assesses the output action along multiple dimensions, namely format, action type, coordinate, and content, to provide detailed feedback. It consists of four components:



**Figure 4** Action type distribution in our mobile training data, showing a long-tailed profile with several low-frequency (sparse) actions.

**Format Reward.** We apply a format reward to train the model to generate reasoning and action outputs that follow a predefined template. This design is motivated by the fact that the reasoning process provides useful historical context for subsequent decision-making. The format reward, denoted as  $R_{\text{format}}$ , checks whether the model output contains the required XML-style tags in the correct order, with the reasoning block enclosed in `<think>` followed by the action block enclosed in `<action>` [DeepSeek-AI \(2025\)](#); [Qian et al. \(2025\)](#).

**Action Type Reward.** We adopt the action type reward,  $R_{\text{type}}$ , which compares the predicted action type with the ground-truth action type. The model receives a reward of 1 for a match and 0 otherwise. This provides a direct and interpretable metric to encourage accurate action type prediction [Lu et al. \(2025b\)](#).

**Coordinate Reward.** For actions involving spatial positioning (*e.g.*, Click, Scroll), we adopt a coordinate reward,  $R_{\text{coord}}$ , that depends on the pixel-level distance between the predicted and ground-truth coordinates. Prior methods often assess correctness by checking whether this distance is within a fixed threshold, which can be somewhat rigid and insufficient for distinguishing near-misses from large deviations. To address this, we introduce a stepwise reward strategy that grants higher scores for predictions closer to the target, while still awarding partial credit for reasonably accurate locations. This design encourages the model to incrementally improve its coordinate parameter localization.

For point-targeting actions (*i.e.*, Click, LongPress), we adopt a stepwise coordinate reward based on the pixel-level distance  $d$  between the predicted and the ground-truth coordinates:

$$R_{\text{coord}} = \begin{cases} \alpha, & d < \delta_2, \\ 0.5\alpha, & \delta_2 \leq d < \delta_1, \\ 0, & \text{otherwise.} \end{cases} \quad (8)$$

Here,  $\delta_1$  and  $\delta_2$  ( $\delta_2 < \delta_1$ ) are distance thresholds that assign higher scores to more accurate predictions, while still granting partial credit for reasonably close locations.

For scrolling actions, we incorporate both spatial accuracy at the start and end positions as well as the correctness of the scrolling direction. The reward  $R_{\text{scroll}}$  is defined as:

$$R_{\text{scroll}} = \begin{cases} 1.5\beta & \text{if } d_{\text{start}}, d_{\text{end}} < \delta_3 \text{ and direction match,} \\ \beta & \text{if } d_{\text{start}} < \delta_3 \text{ and direction match,} \\ 0.5\beta & \text{if } d_{\text{start}} < \delta_3 \text{ or direction match,} \\ 0 & \text{otherwise.} \end{cases} \quad (9)$$

Here,  $d_{\text{start}}$  and  $d_{\text{end}}$  denote the pixel-level distances between the predicted and ground-truth start and end coordinates, respectively, and  $\delta_3$  is the spatial distance threshold for scroll evaluation. The direction refers to the scrolling orientation (*i.e.*, up, down, left, or right), and a direction match indicates that the predicted orientation exactly matches the ground-truth.

**Content Reward.** For actions involving text input (*e.g.*, Type), we compute the reward  $R_{\text{content}}$  using the token-level F1-score between the predicted and ground-truth text:

$$R_{\text{content}} = \begin{cases} \gamma & \text{if F1-score} \geq 0.5, \\ 0 & \text{otherwise.} \end{cases} \quad (10)$$

Here, the F1 score captures the precision-recall balance of textual overlap between the prediction and the reference.

**Total Reward.** Combining all components, the final action-wise reward is computed as:

$$R = R_{\text{format}} \cdot w_1 + (R_{\text{type}} + R_{\text{coord}} + R_{\text{content}}) \cdot w_2, \quad (11)$$

where  $w_1$  and  $w_2$  control the relative importance of structural correctness and action parameter precision.

## 4 Experiments

### 4.1 Implementation Details

We summarize the training hyperparameters in Table 4 and provide the implementation details in the following section.

Hyperparameter	rollout	batch_size	$\beta_{KL}$	freeze_vit	epoch
<b>UI-Venus-Ground-7B</b>	8	128	4e-3	False	10
<b>UI-Venus-Ground-72B</b>	10	128	4e-3	False	1
<b>UI-Venus-Navi-7B</b>	8	256	1e-3	True	10
<b>UI-Venus-Navi-72B</b>	8	512	1e-3	True	1

**Table 4** Training hyperparameter settings used in the experiments of UI-Venus.

**Grounding.** For UI grounding, we train our model with the above-mentioned clean data based on Qwen2.5-VL [Bai et al. \(2025\)](#). Specifically, we set the learning rate as  $4 \times 10^{-7}$  and the rollout sample 8, 10 for 7B, 72B model, respectively. The global batch size is 128 and the models are trained with 128 ppu-gpus with the EasyR1 framework [Zheng et al. \(2025\)](#). It takes 1.5 days, 10 days to train a 7B, 72B model for about one epoch, respectively. To fully reproduce our benchmark results, we provide our evaluation prompt for grounding in Appendix A.1.

Models	Mobile		Desktop		Web		Avg
	Text	Icon/Widget	Text	Icon/Widget	Text	Icon/Widget	
<i>Closed-source Models</i>							
GPT-4o ( <a href="#">OpenAI, 2024</a> )	26.6	24.2	24.2	19.3	12.8	11.8	20.1
UI-TARS-1.5 ( <a href="#">Seed, 2025b</a> )	-	-	-	-	-	-	94.2
Seed1.5-VL ( <a href="#">Seed, 2025a</a> )	-	-	-	-	-	-	<u>95.2</u>
<i>General Open-source Models</i>							
Qwen2.5-VL-7B* ( <a href="#">Bai et al., 2025</a> )	98.3	85.3	88.7	58.6	92.7	81.8	87.7
Qwen2.5-VL-72B* ( <a href="#">Bai et al., 2025</a> )	97.6	88.6	92.3	86.6	91.9	85.2	90.7
<i>GUI-specific Models (SFT)</i>							
SeeClick-9.6B ( <a href="#">Cheng et al., 2024</a> )	78.4	50.7	70.1	29.3	55.2	32.5	55.1
ShowUI-2B ( <a href="#">Lin et al. (2024)</a> )	92.1	75.4	78.9	78.9	84.2	61.1	77.3
UGround-7B ( <a href="#">Gou et al., 2024</a> )	75.1	84.5	85.1	61.4	84.6	71.9	76.3
OS-Atlas-7B ( <a href="#">Wu et al., 2024</a> )	95.2	75.8	90.7	63.6	90.6	77.3	84.1
Aguvis-7B ( <a href="#">Xu et al., 2024</a> )	89.3	68.7	80.6	67.9	89.3	70.0	80.5
UI-TARS-7B ( <a href="#">Qin et al., 2025</a> )	96.9	89.1	95.4	85.0	93.6	85.2	91.6
UI-TARS-72B ( <a href="#">Qin et al., 2025</a> )	94.8	86.3	91.2	87.9	91.5	87.7	90.3
JEDI-7B ( <a href="#">Xie et al., 2025</a> )	96.9	87.2	95.9	87.9	94.4	84.2	91.7
GUI-Actor-7B ( <a href="#">Wu et al., 2025b</a> )	97.6	88.2	96.9	85.7	93.2	86.7	92.1
OpenCUA-7B ( <a href="#">Wang et al., 2025a</a> )	-	-	-	-	-	-	92.3
OpenCUA-32B ( <a href="#">Wang et al., 2025a</a> )	-	-	-	-	-	-	<u>93.4</u>
<i>GUI-specific Models (RL)</i>							
UI-R1-E-3B ( <a href="#">Lu et al., 2025b</a> )	98.2	83.9	94.8	75.0	93.2	83.7	89.5
SE-GUI-7B ( <a href="#">Yuan et al., 2025</a> )	-	-	-	-	-	-	90.3
LPO ( <a href="#">Tang et al., 2025c</a> )	97.9	82.9	95.9	86.4	<u>95.6</u>	84.2	90.5
GUI-G <sup>2</sup> -7B ( <a href="#">Tang et al., 2025a</a> )	-	-	-	-	-	-	93.3
Phi-Ground-7B-16C-DPO ( <a href="#">Zhang et al., 2025a</a> )	96.5	62.0	90.2	76.4	93.6	75.9	83.8
GTA1-7B† ( <a href="#">Yang et al., 2025</a> )	99.0	88.6	94.9	89.3	92.3	86.7	92.4
GTA1-72B ( <a href="#">Yang et al., 2025</a> )	<u>99.3</u>	<u>92.4</u>	<b>97.4</b>	89.3	95.3	<u>91.4</u>	94.8
<i>Ours</i>							
UI-Venus-Ground-7B	99.0	90.0	97.0	<b>90.7</b>	<b>96.2</b>	88.7	94.1
UI-Venus-Ground-72B	<b>99.7</b>	<b>93.8</b>	95.9	<u>90.0</u>	<b>96.2</b>	<b>92.6</b>	<b>95.3</b>

**Table 5** Performance comparison on **ScreenSpot-V2** dataset. Our UI-Venus-72B achieves state-of-the-art performance, outperforming all baseline methods across mobile, desktop, and web platforms. Note that models with \* are reproduced and the † means trained from UI-TARS-1.5-7B.

**Navigation.** Based on the filtered navigation training data, the training is conducted with a learning rate of  $4 \times 10^{-7}$  and a rollout sample size of 8. The global batch size is set to 256 and 512 for the 7B and 72B models, respectively, with training conducted on 256 and 512 PPU-GPUs. Training one epoch takes approximately 1 day for the 7B model and 8.5 days for the 72B model. In addition, we provide the user prompt for the navigation task in Appendix A.2, where the model is constrained to output its reasoning and conclusions along with the predicted actions in a specified format.

## 4.2 Grounding Benchmarks

We evaluate UI-Venus on five comprehensive GUI grounding benchmarks to assess its ability to associate natural language instructions with corresponding GUI elements, including ScreenSpot-V2 [Wu et al. \(2024\)](#), ScreenSpot-Pro [Li et al. \(2025b\)](#), OSWorld-G [Xie et al. \(2025\)](#), UI-Vision [Nayak et al. \(2025\)](#) and CA-GUI [Zhang et al. \(2025b\)](#). During the evaluation, we follow the standard protocol [Cheng et al. \(2024\)](#); [Lin et al. \(2024\)](#), *i.e.*, a prediction is considered correct when the center of predicted box falls within the ground truth bounding box.

In the experiments, we compare UI-Venus models against various state-of-the-art baselines across different model categories: **(1) Closed-source Models:** UI-TARS-1.5 [Seed \(2025b\)](#), Seed1.5-VL [Seed \(2025a\)](#), GPT-4o [OpenAI \(2024\)](#), and Claude Computer Use [Anthropic \(2024\)](#). **(2)**

Model	CAD		Dev		Creative		Scientific		Office		OS		Avg.		
	Text	Icon	Avg.												
<i>Closed-source Models</i>															
GPT-4o ( <a href="#">OpenAI, 2024</a> )	2.0	0.0	1.3	0.0	1.0	0.0	2.1	0.0	1.1	0.0	0.0	0.0	1.3	0.0	0.8
Claude Computer Use ( <a href="#">Anthropic, 2024</a> )	14.5	3.7	22.0	3.9	25.9	3.4	33.9	15.8	30.1	16.3	11.0	4.5	23.4	7.1	17.1
UI-TARS-1.5 ( <a href="#">Seed, 2025b</a> )	-	-	-	-	-	-	-	-	-	-	-	-	-	-	61.6
Seed1.5-VL ( <a href="#">Seed, 2025a</a> )	-	-	-	-	-	-	-	-	-	-	-	-	-	-	60.9
<i>General Open-source Models</i>															
Qwen2.5-VL-7B* ( <a href="#">Bai et al., 2025</a> )	16.8	1.6	46.8	4.1	35.9	7.7	49.3	7.3	52.5	20.8	37.4	6.7	38.9	7.1	26.8
Qwen2.5-VL-72B* ( <a href="#">Bai et al., 2025</a> )	54.8	15.6	65.6	16.6	63.1	19.6	78.5	34.5	79.1	47.2	66.4	29.2	67.3	25.0	51.2
<i>GUI-specific Models (SFT)</i>															
SeeClick-9.6B ( <a href="#">Cheng et al., 2024</a> )	2.5	0.0	0.6	0.0	1.0	0.0	3.5	0.0	1.1	0.0	2.8	0.0	1.8	0.0	1.1
FOCUS-2B ( <a href="#">Tang et al., 2025b</a> )	7.6	3.1	22.8	1.7	23.7	1.7	25.0	7.1	23.2	7.7	17.8	2.5	19.8	3.9	13.3
CogAgent-18B ( <a href="#">Hong et al., 2024</a> )	7.1	3.1	14.9	0.7	9.6	0.0	22.2	1.8	13.0	0.0	5.6	0.0	12.0	0.8	7.7
Aria-UI ( <a href="#">Yang et al., 2024</a> )	7.6	1.6	16.2	0.0	23.7	2.1	27.1	6.4	20.3	1.9	4.7	0.0	17.1	2.0	11.3
OS-Atlas-7B ( <a href="#">Wu et al., 2024</a> )	12.2	4.7	33.1	1.4	28.8	2.8	37.5	7.3	33.9	5.7	27.1	4.5	28.1	4.0	18.9
ShowUI-2B ( <a href="#">Lin et al., 2024</a> )	2.5	0.0	16.9	1.4	9.1	0.0	13.2	7.3	15.3	7.5	10.3	2.2	10.8	2.6	7.7
UGround-7B ( <a href="#">Gou et al., 2024</a> )	14.2	1.6	26.6	2.1	27.3	2.8	31.9	2.7	31.6	11.3	17.8	0.0	25.0	2.8	16.5
UGround-V1-7B ( <a href="#">Gou et al., 2024</a> )	15.8	1.2	51.9	2.8	47.5	9.7	57.6	14.5	60.5	13.2	38.3	7.9	45.2	8.1	31.1
UI-TARS-7B ( <a href="#">Qin et al., 2025</a> )	20.8	9.4	58.4	12.4	50.0	9.1	63.9	31.8	63.3	20.8	30.8	16.9	47.8	16.2	35.7
UI-TARS-72B ( <a href="#">Qin et al., 2025</a> )	18.8	12.5	62.9	17.2	57.1	15.4	64.6	20.9	63.3	26.4	42.1	15.7	50.9	17.6	38.1
JEDI-7B ( <a href="#">Xie et al., 2025</a> )	38.0	14.1	42.9	11.0	50.0	11.9	72.9	25.5	75.1	47.2	33.6	16.9	52.6	18.2	39.5
GUI-Actor-7B ( <a href="#">Wu et al., 2025b</a> )	-	-	-	-	-	-	-	-	-	-	-	-	-	-	44.6
OpenCUA-7B ( <a href="#">Wang et al., 2025a</a> )	-	-	-	-	-	-	-	-	-	-	-	-	-	-	50.0
OpenCUA-32B ( <a href="#">Wang et al., 2025a</a> )	-	-	-	-	-	-	-	-	-	-	-	-	-	-	55.3
<i>GUI-specific Models (RL)</i>															
UI-R1-E-3B ( <a href="#">Lu et al., 2025b</a> )	37.1	12.5	46.1	6.9	41.9	4.2	56.9	21.8	65.0	26.4	32.7	10.1	-	-	33.5
GUI-R1-7B ( <a href="#">Luo et al., 2025a</a> )	23.9	6.3	49.4	4.8	38.9	8.4	55.6	11.8	58.7	26.4	42.1	16.9	-	-	-
InfigUI-R1-3B ( <a href="#">Liu et al., 2025</a> )	33.0	14.1	51.3	12.4	44.9	7.0	58.3	20.0	65.5	28.3	43.9	12.4	49.1	14.1	35.7
GUI-G1-3B ( <a href="#">Zhou et al., 2025</a> )	39.6	9.4	50.7	10.3	36.6	11.9	61.8	30.0	67.2	32.1	23.5	10.6	49.5	16.8	37.1
SE-GUI-7B ( <a href="#">Yuan et al., 2025</a> )	51.3	<b>42.2</b>	68.2	19.3	57.6	9.1	75.0	28.2	78.5	43.4	49.5	25.8	63.5	21.0	47.3
Phi-Ground-7B-16C-DPO ( <a href="#">Zhang et al., 2025a</a> )	26.9	17.2	70.8	16.7	56.6	13.3	58.0	29.1	76.4	44.0	55.1	25.8	56.4	21.8	43.2
GUI-G <sup>2</sup> -7B ( <a href="#">Tang et al., 2025a</a> )	55.8	12.5	68.8	17.2	57.1	15.4	77.1	24.5	74.0	32.7	57.9	21.3	64.7	19.6	47.5
UI-TARS-1.5-7B ( <a href="#">Seed, 2025b</a> )	-	-	-	-	-	-	-	-	-	-	-	-	-	-	49.6
GTA1-7B† ( <a href="#">Yang et al., 2025</a> )	53.3	17.2	66.9	<b>20.7</b>	<b>62.6</b>	18.2	76.4	31.8	82.5	<b>50.9</b>	48.6	25.9	65.5	25.2	50.1
GTA1-72B ( <a href="#">Yang et al., 2025</a> )	56.9	28.1	<b>79.9</b>	<b>33.1</b>	<b>73.2</b>	20.3	81.9	38.2	<b>85.3</b>	49.1	73.8	<b>39.1</b>	<b>74.5</b>	<b>32.5</b>	58.4
<i>Ours</i>															
UI-Venus-Ground-7B	60.4	21.9	74.7	24.1	63.1	14.7	76.4	31.8	75.7	41.5	49.5	22.5	67.1	24.3	50.8
UI-Venus-Ground-72B	<b>66.5</b>	<b>29.7</b>	<b>84.4</b>	<b>33.1</b>	<b>73.2</b>	<b>30.8</b>	<b>84.7</b>	<b>42.7</b>	<b>83.1</b>	<b>60.4</b>	<b>75.7</b>	<b>36.0</b>	<b>77.4</b>	<b>36.8</b>	<b>61.9</b>

**Table 6** Performance comparison of different agent models across various task categories based on Text, Icon, and Average scores on **ScreenSpot-Pro**. Note that models with \* are reproduced and the † means trained from UI-TARS-1.5-7B.

**General Open-source Models:** Qwen2.5-VL-7B/72B [Bai et al. \(2025\)](#) and InternVL2.5 [Chen et al. \(2024\)](#). **(3) GUI-specific SFT Models:** UI-TARS-7B/72B [Qin et al. \(2025\)](#), OS-Atlas-7B [Wu et al. \(2024\)](#), Jedi-7B [Xie et al. \(2025\)](#), and GUI-Actor-7B [Wu et al. \(2025b\)](#). **(4) GUI RL Models:** GTA1-7B/72B [Yang et al. \(2025\)](#), SE-GUI-7B [Yuan et al. \(2025\)](#), and UI-TARS-1.5-7B [Seed \(2025b\)](#). Methods marked with \* indicate our own reproductions, some of which achieve superior performance compared to their original reported results.

**ScreenSpot-V2.** As a general GUI grounding benchmark across mobile, web and desktop platforms with text and icon/widget elements, this benchmark tests the basic grounding capabilities in daily GUI-related scenarios. From Table 5, it can be observed that though previous models have already obtained high scores, UI-Venus-Ground-72B further pushes the performance boundary with 95.3% average score and exhibits superior grounding ability in diverse scenarios of daily life. Notably, UI-Venus-Ground-7B, with significantly fewer parameters, can still achieve 94.1% and outperform larger models like UI-TARS-72B or OpenCUA-32B, which only score 90.3% and 93.4%, respectively.

Models	Text Matching	Element Recognition	Layout Understanding	Fine-grained Manipulation	Refusal	Avg
<i>Closed-source Models</i>						
Operator (OpenAI, 2025)	51.3	42.4	46.6	31.5	-	40.6
Gemini-2.5-Pro (Team, 2024)	59.8	45.5	49.0	33.6	<b>38.9</b>	45.2
Seed1.5-VL (Seed, 2025a)	73.9	66.7	69.6	47.0	<u>18.5</u>	62.9
<i>Open-source Models</i>						
OS-Atlas-7B (Wu et al., 2024)	44.1	29.4	35.2	16.8	7.4	27.7
Qwen2.5-VL-7B (Bai et al., 2025)	45.6	32.7	41.9	18.1	-	31.4
Qwen2.5-VL-72B (Bai et al., 2025)	52.6	74.6	<u>74.7</u>	55.3	-	62.2
UGround-7B (Gou et al., 2024)	51.3	40.3	43.5	24.8	-	36.4
Aguvis-7B (Xu et al., 2024)	55.9	41.2	43.9	28.2	-	38.7
Jedi-7B (Xie et al., 2025)	65.9	55.5	57.7	46.9	7.4	54.1
UI-TARS-7B (Qin et al., 2025)	60.2	51.8	54.9	35.6	-	47.5
UI-TARS-1.5-7B (Seed, 2025b)	52.6	75.4	72.4	<u>66.7</u>	-	64.2
UI-TARS-72B (Qin et al., 2025)	69.4	60.6	62.9	<u>45.6</u>	-	57.1
GTA1-7B (Yang et al., 2025)	63.2	<b>82.1</b>	74.2	<b>70.5</b>	-	<u>67.7</u>
GTA1-72B (Yang et al., 2025)	57.9	<u>76.9</u>	<b>77.3</b>	<u>66.7</u>	-	66.7
OpenCUA-7B (Wang et al., 2025a)	-	-	-	-	-	55.3
OpenCUA-32B (Wang et al., 2025a)	-	-	-	-	-	59.6
<i>Ours</i>						
UI-Venus-Ground-7B	<u>74.6</u>	60.5	61.5	45.5	-	58.8
UI-Venus-Ground-72B	<b>82.1</b>	71.2	70.7	64.4	-	<b>70.4</b>

**Table 7** Performance comparison on **OS-World-G** grounding dataset. Our UI-Venus-72B achieves state-of-the-art performance, outperforming all baseline methods across different settings.

**ScreenSpot-Pro.** With high-resolution professional software interfaces, including CAD, development, creative, scientific, office, and OS applications, the difficulty of element grounding is greatly increased. Nevertheless, as shown in Table 6, UI-Venus-Ground-72B can achieve comprehensive leadership with 61.9% score. Compared to text elements, we can find in professional software, icon elements are much harder to detect and locate due to their diverse and small shapes, which requires the model’s understanding of GUI layouts and fine-grained coordinate perception, while UI-Venus-Ground-72B still makes a breakthrough in icon grounding for software of ‘Creative’, ‘Scientific’ and ‘Office’ categories, with improvements over previous optimal model GTA1-72B of 10.5%, 4.5% and 9.5%, respectively.

**OSWorld-G.** This benchmark is composed of fine-grained tasks from real computer environment with diverse capabilities including text matching, element recognition, layout understanding, fine-grained manipulation and refusal handling. The result in Table 7 demonstrates our model’s exceptional performance in such tasks, where 72B model achieves 70.4% score, substantially surpassing strong baselines like UI-TARS-1.5-7B and GTA1 series.

**UI-Vision.** UI-Vision provides a fine-grained evaluation of computer-using agents in real-world desktop environments, serving as a comprehensive, license-permissive benchmark. In Table 8, we can observe UI-Venus-Ground-72B establishes new state-of-the-art performance across three task categories, while UI-Venus-Ground-7B performs on par with previous best model Phi-Ground-7B-16C-DPO, indicating our model has made significant progress in grounding capability under real-world environments.

**CA-GUI.** We also adopt multilingual evaluation with realistic Chinese mobile applications, featuring Fun2Point and Text2Point tasks, and the results are exhibited in Table 9. In this benchmark, our model demonstrates strong cross-lingual generalization capabilities in non-English GUI environments, markedly exceeding AgentCPM-GUI by 5% in Fun2Point task and by 9.4% in Text2Point

Models	Basic	Functional	Spatial	Avg
<i>Closed-source Models</i>				
GPT-4o ( <a href="#">OpenAI, 2024</a> )	1.6	1.5	1.0	1.38
Claude-3.7-Sonnet ( <a href="#">Anthropic, 2024</a> )	9.5	7.7	7.6	8.3
<i>Open-source Models</i>				
Qwen2.5-VL-7B ( <a href="#">Bai et al., 2025</a> )	1.2	0.8	0.5	0.9
SeeClick ( <a href="#">Cheng et al., 2024</a> )	9.4	4.7	2.1	5.4
UGround-V1-7B ( <a href="#">Gou et al., 2024</a> )	15.4	17.1	6.3	12.9
OS-Atlas-7B ( <a href="#">Wu et al., 2024</a> )	12.2	11.2	3.7	9.0
UI-TARS-7B ( <a href="#">Qin et al., 2025</a> )	20.1	24.3	8.4	17.6
UI-TARS-1.5-7B ( <a href="#">Seed, 2025b</a> )	28.8	27.5	10.7	22.3
UI-TARS-72B ( <a href="#">Qin et al., 2025</a> )	31.4	30.5	14.7	25.5
Phi-Ground-7B-16C-DPO ( <a href="#">Zhang et al., 2025a</a> )	36.8	37.1	7.6	27.2
<i>Ours</i>				
UI-Venus-Ground-7B	36.1	32.8	11.9	26.5
UI-Venus-Ground-72B	<b>45.6</b>	<b>42.3</b>	<b>23.7</b>	<b>36.8</b>

**Table 8** Performance comparison on **UI-Vision** grounding dataset. Our UI-Venus-Ground-72B achieves state-of-the-art performance, outperforming all baseline methods across different settings.

Models	Fun2Point	Text2Point	Avg
<i>Closed-source Models</i>			
GPT-4o ( <a href="#">OpenAI, 2024</a> )	22.1	19.9	21.0
GPT-4o w grounding ( <a href="#">OpenAI, 2024</a> )	44.3	44.0	44.2
<i>Open-source Models</i>			
Intern2.5-VL-8B ( <a href="#">Chen et al., 2024</a> )	17.2	24.2	20.7
Intern2.5-VL-26B ( <a href="#">Chen et al., 2024</a> )	14.8	16.6	15.7
Qwen2.5-VL-7B ( <a href="#">Bai et al., 2025</a> )	59.8	59.3	59.6
OS-Genesis-7B ( <a href="#">Sun et al., 2025a</a> )	8.3	5.8	7.1
OS-Atlas-7B ( <a href="#">Wu et al., 2024</a> )	53.6	60.7	57.2
Aguvis-7B ( <a href="#">Xu et al., 2024</a> )	60.8	76.5	68.7
UI-TARS-7B ( <a href="#">Qin et al., 2025</a> )	56.8	66.7	61.8
AgentCPM-GUI ( <a href="#">Zhang et al., 2025b</a> )	79.1	76.5	77.8
<i>Ours</i>			
UI-Venus-Ground-7B	83.3	83.2	83.3
UI-Venus-Ground-72B	<b>84.1</b>	<b>85.9</b>	<b>85.0</b>

**Table 9** Performance comparison on **CA-GUI** grounding dataset. Our UI-Venus-72B achieves state-of-the-art performance, outperforming all baseline methods across different settings.

tasks, and even UI-Venus-Ground-7B shows great advantages.

In summary, our evaluation reveals several key insights: **(1) Consistent SOTA Performance:** UI-Venus achieves new state-of-the-art results across all benchmarks, with UI-Venus-Ground-72B reaching 95.3% on ScreenSpot-V2, 61.9% on ScreenSpot-Pro, and 85.0% on CA-GUI, significantly outperforming previous best models including GTA1-72B (94.8%, 58.4%, -) and UI-TARS-1.5 (95.2%, 61.6%, -). **(2) Superior Fine-grained Capabilities:** On OSWorld-G, our model demonstrates exceptional performance in fine-grained manipulation tasks, substantially surpassing strong baselines like UI-TARS-1.5-7B and GTA1 series. **(3) Robust Multilingual Grounding:** UI-Venus shows substantial improvements on the Chinese benchmark CA-GUI, achieving 83.2% average accuracy and demonstrating strong cross-lingual generalization capabilities in non-English GUI environments.

Models	With Planner	A11y Tree	Screenshot	Success Rate(pass@1)
<i>Closed-source Models</i>				
GPT-4o ( <a href="#">OpenAI, 2024</a> )	✗	✓	✗	30.6
ScaleTrack ( <a href="#">Huang et al., 2025</a> )	✗	✓	✗	44.0
SeedVL-1.5 ( <a href="#">Seed, 2025a</a> )	✗	✓	✓	62.1
UI-TARS-1.5 ( <a href="#">Seed, 2025b</a> )	✗	✗	✓	64.2
<i>Open-source Models</i>				
GUI-Critic-R1-7B ( <a href="#">Wanyan et al., 2025</a> )	✗	✓	✓	27.6
Qwen2.5-VL-72B ( <a href="#">Bai et al., 2025</a> )	✗	✗	✓	35.0
UGround ( <a href="#">Gou et al., 2024</a> )	✓	✗	✓	44.0
Aria-UI ( <a href="#">Yang et al., 2024</a> )	✓	✗	✓	44.8
UI-TARS-72B ( <a href="#">Qin et al., 2025</a> )	✗	✗	✓	46.6
GLM-4.5v ( <a href="#">Zhipu-AI, 2025</a> )	✗	✗	✓	57.0
<i>Ours</i>				
UI-Venus-Navi-7B	✗	✗	✓	49.1
UI-Venus-Navi-72B	✗	✗	✓	<b>65.9</b>

**Table 10** Performance comparison on **AndroidWorld** for end-to-end models. Our UI-Venus-Navi-72B achieves state-of-the-art performance, outperforming all baseline methods across different settings.

### 4.3 Navigation Benchmarks

We evaluate UI-Venus on three widely-adopted benchmarks: AndroidWorld [Rawles et al. \(2025\)](#), AndroidControl [Li et al. \(2024\)](#), and GUI-Odyssey [Lu et al. \(2025a\)](#) datasets to assess its multi-step decision-making capabilities across diverse mobile interface scenarios. Methods marked with \* indicate our own reproductions, some of which achieve superior performance compared to their original reported results.

**Online Benchmark.** For real-time multi-step decision-making evaluation, we adopt AndroidWorld, a dynamic benchmark requiring continuous interaction with live mobile applications. This framework assesses the model’s ability to adapt strategies based on real-time feedback and maintain task coherence across extended interaction sequences. From Table 10, we can see UI-Venus is implemented without relying on additional planner or A11y tree and is able to conduct end-to-end UI navigation purely based on screenshots, which ensures strong generalization capability in real-time interactive scenarios. In AndroidWorld benchmark, UI-Venus-Navi-72B gains a score of 65.9% success rate in one-time evaluation, surpassing UI-TARS-1.5 (64.2%), and UI-Venus-Navi-7B also succeeds with 49.1% rate to outperform UI-TARS-72B (46.6%), demonstrating the effectiveness of UI-Venus.

**Offline Benchmark.** Besides, we employ two established offline benchmarks: AndroidControl and GUI-Odyssey. These static evaluation scenarios assess the model’s fundamental UI understanding, task decomposition, and action planning capabilities in controlled environments. With evaluation results in Table 11, where ‘Low/High’ categories in AndroidControl represent ‘the low/high level instruction’ the agent is prompted with, respectively, we can observe that UI-Venus produces comparable results with previous optimal methods. Notably, in AndroidControl-High evaluation, UI-Venus performs best in both type accuracy and step success rate, indicating that UI-Venus possesses stronger planning and summary capability with high-level task instruction.

In summary, our evaluation reveals several key insights: (1) **SOTA Online Navigation Performance:** UI-Venus achieves new state-of-the-art end-to-end model results on AndroidWorld, with UI-Venus-Navi-72B reaching **65.9%** success rate, significantly outperforming previous best models including

Models	AndroidControl-Low		AndroidControl-High		GUI-Odyssey	
	Type Acc.	Step SR	Type Acc.	Step SR	Type Acc.	Step SR
<i>Closed-source Models</i>						
GPT-4o (OpenAI, 2024)	74.3	19.4	66.3	20.8	34.3	3.3
<i>Open-source Models</i>						
Qwen2.5-VL-7B (Bai et al., 2025)	94.1	85.0	75.1	62.9	59.5	46.3
SeeClick (Cheng et al., 2024)	93.0	75.0	82.9	59.1	71.0	53.9
OS-Atlas-7B (Wu et al., 2024)	93.6	85.2	85.2	71.2	84.5	62.0
Aguvis-7B (Xu et al., 2024)	-	80.5	-	61.5	-	-
Aguvis-72B (Xu et al., 2024)	-	84.4	-	66.4	-	-
OS-Genesis-7B (Xu et al., 2024)	90.7	74.2	66.2	44.5	-	-
UI-TARS-7B (Qin et al., 2025)	<u>98.0</u>	90.8	83.7	72.5	<u>94.6</u>	87.0
UI-TARS-72B (Qin et al., 2025)	<b>98.1</b>	91.3	85.2	74.7	<b>95.4</b>	88.6
GUI-R1-7B (Luo et al., 2025a)	85.2	66.5	71.6	51.7	65.5	38.8
NaviMaster-7B (Luo et al., 2025b)	85.6	69.9	72.9	54.0	-	-
UI-AGILE-7B (Lian et al., 2025)	87.7	77.6	80.1	60.6	-	-
AgentCPM-GUI (Zhang et al., 2025b)	94.4	90.2	77.7	69.2	90.0	<b>75.0</b>
<i>Ours</i>						
UI-Venus-Navi-7B	97.1	<u>92.4</u>	<b>86.5</b>	<u>76.1</u>	87.3	71.5
UI-Venus-Navi-72B	96.7	<b>92.9</b>	<u>85.9</u>	<u>77.2</u>	87.2	<u>72.4</u>

**Table 11** Performance comparison on offline UI navigation datasets including **AndroidControl** and **GUI-Odyssey**.

UI-TARS-72B (46.6%) and UI-TARS-1.5 (64.2%). We also show some case studies in Sec. B.2 and open-source the all evaluation traces of AndroidWorld to better show the effectiveness of our UI-Venus. **(2) Comparable Offline Navigation Performance:** On AndroidControl and GUI-Odyssey, our UI-Venus also achieves comparable results with previous SOTA methods. Note that UI-Venus obtains best performance on AndroidControl-High, reflecting its remarkable generalization performance of planning and summary in long traces.

## 5 Future Work

Although the GRPO algorithm has proven to be a more effective post-training strategy for UI agents than conventional supervised fine-tuning, certain challenges remain. One notable issue is the hallucination gap between the model’s internal reasoning (*think*) and its final response (*answer*), which can lead to incorrect or inconsistent behavior in navigation tasks. Another observation is that even humans may struggle to operate unfamiliar applications without prior exposure. This suggests that large-scale pretraining on trajectories from diverse UI agents can equip models with richer prior knowledge of various applications and to enhance their adaptability. Thus, future research may explore integrating pretraining with curated UI-agent traces, refining action generation through enhanced reasoning alignment, and incorporating domain-specific priors. Such advancements would help mitigate these issues and enable broader deployment of UI agents across diverse and dynamic computing environments.

## References

- Droidrun. <https://github.com/droidrun/droidrun>, 2025.
- Saket Agashe, Jiuzhou Han, Shuyu Gan, Jiachen Yang, Ang Li, and Xin Eric Wang. Agent s: An open agentic framework that uses computers like a human. *arXiv preprint arXiv:2410.08164*, 2024.
- Anthropic. Claude computer use. Available at: <https://www.anthropic.com/news/developing-computer-use>, 2024.
- Chongyang Bai, Xiaoxue Zang, Ying Xu, Srinivas Sunkara, Abhinav Rastogi, Jindong Chen, and Blaise Aguera y Arcas. Ubert: Learning generic multimodal representations for ui understanding, 2021. <https://arxiv.org/abs/2107.13731>.
- Hao Bai, Yifei Zhou, Mert Cemri, Jiayi Pan, Alane Suhr, Sergey Levine, and Aviral Kumar. Digirl: Training in-the-wild device-control agents with autonomous reinforcement learning, 2024. <https://arxiv.org/abs/2406.11896>.
- 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, 2023. <https://arxiv.org/abs/2308.12966>.
- 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, 2025. <https://arxiv.org/abs/2502.13923>.
- Hyungjoo Chae, Namyoung Kim, Kai Tzu iunn Ong, Minju Gwak, Gwanwoo Song, Jihoon Kim, Sunghwan Kim, Dongha Lee, and Jinyoung Yeo. Web agents with world models: Learning and leveraging environment dynamics in web navigation. In *The Thirteenth International Conference on Learning Representations*, 2025. <https://openreview.net/forum?id=moWiYJuSGF>.
- Hardy Chen, Haoqin Tu, Fali Wang, Hui Liu, Xianfeng Tang, Xinya Du, Yuyin Zhou, and Cihang Xie. Sft or rl? an early investigation into training rl-like reasoning large vision-language models, 2025. <https://arxiv.org/abs/2504.11468>.
- Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024.
- Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Yantao Li, Jianbing Zhang, and Zhiyong Wu. Seeclick: Harnessing gui grounding for advanced visual gui agents, 2024. <https://arxiv.org/abs/2401.10935>.
- DeepSeek-AI. Deepseek-rl: Incentivizing reasoning capability in llms via reinforcement learning, 2025. <https://arxiv.org/abs/2501.12948>.
- Lang Feng, Weihao Tan, Zhiyi Lyu, Longtao Zheng, Haiyang Xu, Ming Yan, Fei Huang, and Bo An. Towards efficient online tuning of vlm agents via counterfactual soft reinforcement learning, 2025. <https://arxiv.org/abs/2505.03792>.
- Minghe Gao, Wendong Bu, Bingchen Miao, Yang Wu, Yunfei Li, Juncheng Li, Siliang Tang, Qi Wu, Yueting Zhuang, and Meng Wang. Generalist virtual agents: A survey on autonomous agents across digital platforms. *arXiv preprint arXiv:2411.10943*, 2024.
- Boyu Gou, Ruohan Wang, Boyuan Zheng, Yanan Xie, Cheng Chang, Yiheng Shu, Huan Sun, and Yu Su. Navigating the digital world as humans do: Universal visual grounding for gui agents, 2024. <https://arxiv.org/abs/2410.05243>.
- Zhangxuan Gu, Zhuoer Xu, Haoxing Chen, Jun Lan, Changhua Meng, and Weiqiang Wang. Mobile user interface element detection via adaptively prompt tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11155–11164, 2023.
- Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxuan Zhang, Juanzi Li, Bin Xu, Yuxiao Dong, Ming Ding, and Jie Tang. Cogagent: A visual language model for gui agents, 2024. <https://arxiv.org/abs/2312.08914>.
- Xueyu Hu, Tao Xiong, Biao Yi, Zishu Wei, Ruixuan Xiao, Yurun Chen, Jiasheng Ye, Meiling Tao, Xiangxin Zhou, Ziyu Zhao, et al. Os agents: A survey on mllm-based agents for computer, phone and browser use, 2024.
- Jing Huang, Zhixiong Zeng, Wenkang Han, Yufeng Zhong, Liming Zheng, Shuai Fu, Jingyuan Chen, and Lin Ma. Scaletrack: Scaling and back-tracking automated gui agents. *arXiv preprint arXiv:2505.00416*, 2025.
- Peter C Humphreys, David Raposo, Tobias Pohlen, Gregory Thornton, Rachita Chhaparia, Alistair Muldal, Josh Abramson, Petko Georgiev, Adam Santoro, and Timothy Lillicrap. A data-driven approach for learning to control computers. In *International Conference on Machine Learning*, pages 9466–9482. PMLR, 2022.

- Raghav Kapoor, Yash Parag Butala, Melisa Russak, Jing Yu Koh, Kiran Kamble, Waseem Alshikh, and Ruslan Salakhutdinov. Omniact: A dataset and benchmark for enabling multimodal generalist autonomous agents for desktop and web, 2024. <https://arxiv.org/abs/2402.17553>.
- Cheryl Li, Tianyuan Xu, and Yiwen Guo. Reasoning-as-logic-units: Scaling test-time reasoning in large language models through logic unit alignment, 2025a. <https://arxiv.org/abs/2502.07803>.
- Kaixin Li, Ziyang Meng, Hongzhan Lin, Ziyang Luo, Yuchen Tian, Jing Ma, Zhiyong Huang, and Tat-Seng Chua. Screenspot-pro: Gui grounding for professional high-resolution computer use, 2025b.
- Ning Li, Xiangmou Qu, Jiamu Zhou, Jun Wang, Muning Wen, Kounianhua Du, Xingyu Lou, Qiuying Peng, Jun Wang, and Weinan Zhang. Mobileuse: A gui agent with hierarchical reflection for autonomous mobile operation. *arXiv preprint arXiv:2507.16853*, 2025c. <https://arxiv.org/abs/2507.16853>.
- Wei Li, William Bishop, Alice Li, Chris Rawles, Folawiyo Campbell-Ajala, Divya Tyamagundlu, and Oriana Riva. On the effects of data scale on ui control agents, 2024. <https://arxiv.org/abs/2406.03679>.
- Shuquan Lian, Yuhang Wu, Jia Ma, Zihan Song, Bingqi Chen, Xiawu Zheng, and Hui Li. Ui-agile: Advancing gui agents with effective reinforcement learning and precise inference-time grounding, 2025. <https://arxiv.org/abs/2507.22025>.
- 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. <https://arxiv.org/abs/2411.17465>.
- Yuhang Liu, Pengxiang Li, Congkai Xie, Xavier Hu, Xiaotian Han, Shengyu Zhang, Hongxia Yang, and Fei Wu. Infigui-r1: Advancing multimodal gui agents from reactive actors to deliberative reasoners. 2025. <https://arxiv.org/abs/2504.14239>.
- Quanfeng Lu, Wenqi Shao, Zitao Liu, Lingxiao Du, Fanqing Meng, Boxuan Li, Botong Chen, Siyuan Huang, Kaipeng Zhang, and Ping Luo. Guiodyssey: A comprehensive dataset for cross-app gui navigation on mobile devices, 2025a. <https://arxiv.org/abs/2406.08451>.
- Zhengxi Lu, Yuxiang Chai, Yaxuan Guo, Xi Yin, Liang Liu, Hao Wang, Han Xiao, Shuai Ren, Guanjing Xiong, and Hongsheng Li. Ui-r1: Enhancing efficient action prediction of gui agents by reinforcement learning. 2025b. <https://arxiv.org/abs/2503.21620>.
- Run Luo, Lu Wang, Wanwei He, and Xiaobo Xia. Gui-r1 : A generalist r1-style vision-language action model for gui agents. 2025a. <https://arxiv.org/abs/2504.10458>.
- Zhihao Luo, Wentao Yan abd Jingyu Gong, Min Wang, Zhizhong Zhang, Xuhong Wang, Yuan Xie, and Xin Tan. Navimaster: Learning a unified policy for gui and embodied navigation tasks, 2025b. <https://arxiv.org/abs/2508.02046>.
- Shravan Nayak, Xiangru Jian, Kevin Qinghong Lin, Juan A Rodriguez, Montek Kalsi, Rabiul Awal, Nicolas Chapados, M Tamer Özsu, Aishwarya Agrawal, David Vazquez, et al. Ui-vision: A desktop-centric gui benchmark for visual perception and interaction. *arXiv preprint arXiv:2503.15661*, 2025.
- Dang Nguyen, Jian Chen, Yu Wang, Gang Wu, Namyong Park, Zhengmian Hu, Hanjia Lyu, Junda Wu, Ryan Aponte, Yu Xia, et al. Gui agents: A survey. *arXiv preprint arXiv:2412.13501*, 2024.
- OpenAI. Introducing gpt-4o. Available at: <https://openai.com/index/hello-gpt-4o>, 2024.
- OpenAI. Introducing operator. Available at: <https://openai.com/index/introducing-operator/>, 2025.
- Zehan Qi, Xiao Liu, Iat Long Iong, Hanyu Lai, Xueqiao Sun, Wenyi Zhao, Yu Yang, Xinyue Yang, Jiadai Sun, Shuntian Yao, Tianjie Zhang, Wei Xu, Jie Tang, and Yuxiao Dong. Webrl: Training llm web agents via self-evolving online curriculum reinforcement learning, 2025. <https://arxiv.org/abs/2411.02337>.
- Cheng Qian, Emre Can Acikgoz, Qi He, Hongru Wang, Xiusi Chen, Dilek Hakkani-Tür, Gokhan Tur, and Heng Ji. Toolrl: Reward is all tool learning needs, 2025. <https://arxiv.org/abs/2504.13958>.
- Yujia Qin, Yining Ye, Junjie Fang, Haoming Wang, Shihao Liang, Shizuo Tian, Junda Zhang, Jiahao Li, Yunxin Li, Shijue Huang, Wanjun Zhong, Kuanye Li, Jiale Yang, Yu Miao, Woyu Lin, Longxiang Liu, Xu Jiang, Qianli Ma, Jingyu Li, Xiaojun Xiao, Kai Cai, Chuang Li, Yaowei Zheng, Chaolin Jin, Chen Li, Xiao Zhou, Minchao Wang, Haoli Chen, Zhaojian Li, Haihua Yang, Haifeng Liu, Feng Lin, Tao Peng, Xin Liu, and Guang Shi. Ui-tars: Pioneering automated gui interaction with native agents, 2025. <https://arxiv.org/abs/2501.12326>.
- Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy Lillicrap. Androidinthewild: A large-scale dataset for android device control. *Advances in Neural Information Processing Systems*, 36:59708–59728, 2023a.
- Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy Lillicrap. Android in the wild: A large-scale dataset for android device control, 2023b. <https://arxiv.org/abs/2307.10088>.

Christopher Rawles, Sarah Clinckemaillie, Yifan Chang, Jonathan Waltz, Gabrielle Lau, Marybeth Fair, Alice Li, William Bishop, Wei Li, Folawiyo Campbell-Ajala, Daniel Toyama, Robert Berry, Divya Tyamagundlu, Timothy Lillicrap, and Oriana Riva. Androidworld: A dynamic benchmarking environment for autonomous agents, 2025. <https://arxiv.org/abs/2405.14573>.

ByteDance Seed. Seed1.5-vl technical report. *arXiv preprint arXiv:2505.07062*, 2025a.

ByteDance Seed. Ui-tars-1.5. <https://seed-tars.com/1.5>, 2025b.

Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathematical reasoning in open language models, 2024. <https://arxiv.org/abs/2402.03300>.

Haozhan Shen, Peng Liu, Jingcheng Li, Chunxin Fang, Yibo Ma, Jiajia Liao, Qiaoli Shen, Zilun Zhang, Kangjia Zhao, Qianqian Zhang, Ruochen Xu, and Tiancheng Zhao. Vlm-r1: A stable and generalizable r1-style large vision-language model, 2025. <https://arxiv.org/abs/2504.07615>.

Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning, 2023.

Qiushi Sun, Kanzhi Cheng, Zichen Ding, Chuanyang Jin, Yian Wang, Fangzhi Xu, Zhenyu Wu, Chengyou Jia, Liheng Chen, Zhoumianze Liu, Ben Kao, Guohao Li, Junxian He, Yu Qiao, and Zhiyong Wu. Os-genesis: Automating gui agent trajectory construction via reverse task synthesis, 2025a. <https://arxiv.org/abs/2412.19723>.

Qiushi Sun, Kanzhi Cheng, Zichen Ding, Chuanyang Jin, Yian Wang, Fangzhi Xu, Zhenyu Wu, Chengyou Jia, Liheng Chen, Zhoumianze Liu, Ben Kao, Guohao Li, Junxian He, Yu Qiao, and Zhiyong Wu. Os-genesis: Automating gui agent trajectory construction via reverse task synthesis, 2025b. <https://arxiv.org/abs/2412.19723>.

Yuchen Sun, Shanhui Zhao, Tao Yu, Hao Wen, Samith Va, Mengwei Xu, Yuanchun Li, and Chongyang Zhang. Gui-xplore: Empowering generalizable gui agents with one exploration, 2025c. <https://arxiv.org/abs/2503.17709>.

Weihao Tan, Wentao Zhang, Xinrun Xu, Haochong Xia, Ziluo Ding, Boyu Li, Bohan Zhou, Junpeng Yue, Jiechuan Jiang, Yewen Li, et al. Cradle: Empowering foundation agents towards general computer control. *arXiv preprint arXiv:2403.03186*, 2024.

Fei Tang, Zhangxuan Gu, Zhengxi Lu, Xuyang Liu, Shuheng Shen, Changhua Meng, Wen Wang, Wenqi Zhang, Yongliang Shen, Weiming Lu, Jun Xiao, and Yueling Zhuang. Gui-g<sup>2</sup>: Gaussian reward modeling for gui grounding, 2025a. <https://arxiv.org/abs/2507.15846>.

Fei Tang, Yongliang Shen, Hang Zhang, Siqi Chen, Guiyang Hou, Wenqi Zhang, Wenqiao Zhang, Kaitao Song, Weiming Lu, and Yueling Zhuang. Think twice, click once: Enhancing gui grounding via fast and slow systems. 2025b. <https://arxiv.org/abs/2503.06470>.

Jiaqi Tang, Yu Xia, Yi-Feng Wu, Yuwei Hu, Yuhui Chen, Qing-Guo Chen, Xiaogang Xu, Xiangyu Wu, Hao Lu, Yanqing Ma, Shiyin Lu, and Qifeng Chen. Lpo: Towards accurate gui agent interaction via location preference optimization, 2025c. <https://arxiv.org/abs/2506.09373>.

Zhengwei Tao, Ting-En Lin, Xiancai Chen, Hangyu Li, Yuchuan Wu, Yongbin Li, Zhi Jin, Fei Huang, Dacheng Tao, and Jingren Zhou. A survey on self-evolution of large language models, 2024. <https://arxiv.org/abs/2404.14387>.

Gemini Team. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context, 2024. <https://arxiv.org/abs/2403.05530>.

Junyang Wang, Haiyang Xu, Haitao Jia, Xi Zhang, Ming Yan, Weizhou Shen, Ji Zhang, Fei Huang, and Jitao Sang. Mobile-agent-v2: Mobile device operation assistant with effective navigation via multi-agent collaboration, 2024a. <https://arxiv.org/abs/2406.01014>.

Junyang Wang, Haiyang Xu, Jiabo Ye, Ming Yan, Weizhou Shen, Ji Zhang, Fei Huang, and Jitao Sang. Mobile-agent: Autonomous multi-modal mobile device agent with visual perception. *arXiv preprint arXiv:2401.16158*, 2024b.

Ke Wang, Tianyu Xia, Zhangxuan Gu, Yi Zhao, Shuheng Shen, Changhua Meng, Weiqiang Wang, and Ke Xu. E-ant: A large-scale dataset for efficient automatic gui navigation, 2024c. <https://arxiv.org/abs/2406.14250>.

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, 2024d. <https://arxiv.org/abs/2409.12191>.

Shuai Wang, Weiwen Liu, Jingxuan Chen, Yuqi Zhou, Weinan Gan, Xingshan Zeng, Yuhan Che, Shuai Yu, Xinlong Hao, Kun Shao, et al. Gui agents with foundation models: A comprehensive survey. *arXiv preprint arXiv:2411.04890*, 2024e.

- Xinyuan Wang, Bowen Wang, Dunjie Lu, Junlin Yang, Tianbao Xie, Junli Wang, Jiaqi Deng, Xiaole Guo, Yiheng Xu, Chen Henry Wu, Zhennan Shen, Zhuokai Li, Ryan Li, Xiaochuan Li, Junda Chen, Boyuan Zheng, Peihang Li, Fangyu Lei, Ruisheng Cao, Yeqiao Fu, Dongchan Shin, Martin Shin, Jiarui Hu, Yuyan Wang, Jixuan Chen, Yuxiao Ye, Danyang Zhang, Dikang Du, Hao Hu, Huarong Chen, Zaida Zhou, Yipu Wang, Heng Wang, Diyi Yang, Victor Zhong, Flood Sung, Y. Charles, Zhilin Yang, and Tao Yu. Opencua: Open foundations for computer-use agents, 2025a. <https://arxiv.org/abs/2508.09123>.
- Zhenhailong Wang, Haiyang Xu, Junyang Wang, Xi Zhang, Ming Yan, Ji Zhang, Fei Huang, and Heng Ji. Mobile-agent-e: Self-evolving mobile assistant for complex tasks, 2025b. <https://arxiv.org/abs/2501.11733>.
- Yuyang Wanyan, Xi Zhang, Haiyang Xu, Haowei Liu, Junyang Wang, Jiabo Ye, Yutong Kou, Ming Yan, Fei Huang, Xiaoshan Yang, et al. Look before you leap: A gui-critic-r1 model for pre-operative error diagnosis in gui automation. *arXiv preprint arXiv:2506.04614*, 2025.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023. <https://arxiv.org/abs/2201.11903>.
- Penghao Wu, Shengnan Ma, Bo Wang, Jiaheng Yu, Lewei Lu, and Ziwei Liu. Gui-reflection: Empowering multimodal gui models with self-reflection behavior. *arXiv preprint arXiv:2506.08012*, 2025a.
- Qianhui Wu, Kanzhi Cheng, Rui Yang, Chaoyun Zhang, Jianwei Yang, Huiqiang Jiang, Jian Mu, Baolin Peng, Bo Qiao, Reuben Tan, et al. Gui-actor: Coordinate-free visual grounding for gui agents. *arXiv preprint arXiv:2506.03143*, 2025b.
- Zhiyong Wu, Zhenyu Wu, Fangzhi Xu, Yian Wang, Qiushi Sun, Chengyou Jia, Kanzhi Cheng, Zichen Ding, Liheng Chen, Paul Pu Liang, and Yu Qiao. Os-atlas: A foundation action model for generalist gui agents, 2024. <https://arxiv.org/abs/2410.23218>.
- Tianbao Xie, Jiaqi Deng, Xiaochuan Li, Junlin Yang, Haoyuan Wu, Jixuan Chen, Wenjing Hu, Xinyuan Wang, Yuhui Xu, Zekun Wang, Yiheng Xu, Junli Wang, Doyen Sahoo, Tao Yu, and Caiming Xiong. Scaling computer-use grounding via user interface decomposition and synthesis, 2025. <https://arxiv.org/abs/2505.13227>.
- Yiheng Xu, Zekun Wang, Junli Wang, Dunjie Lu, Tianbao Xie, Amrita Saha, Doyen Sahoo, Tao Yu, and Caiming Xiong. Aguvis: Unified pure vision agents for autonomous gui interaction. 2024. <https://arxiv.org/abs/2412.04454>.
- Yan Yang, Dongxu Li, Yutong Dai, Yuhao Yang, Ziyang Luo, Zirui Zhao, Zhiyuan Hu, Junzhe Huang, Amrita Saha, Zeyuan Chen, Ran Xu, Liyuan Pan, Caiming Xiong, and Junnan Li. Gta1: Gui test-time scaling agent, 2025. <https://arxiv.org/abs/2507.05791>.
- Yuhao Yang, Yue Wang, Dongxu Li, Ziyang Luo, Bei Chen, Chao Huang, and Junnan Li. Aria-ui: Visual grounding for gui instructions, 2024. <https://arxiv.org/abs/2412.16256>.
- Xinbin Yuan, Jian Zhang, Kaixin Li, Zhuoxuan Cai, Lujian Yao, Jie Chen, Enguang Wang, Qibin Hou, Jinwei Chen, Peng-Tao Jiang, and Bo Li. Enhancing visual grounding for gui agents via self-evolutionary reinforcement learning. 2025. <https://arxiv.org/abs/2505.12370>.
- Chaoyun Zhang, Shilin He, Jiaxu Qian, Bowen Li, Liquun Li, Si Qin, Yu Kang, Minghua Ma, Guyue Liu, Qingwei Lin, et al. Large language model-brained gui agents: A survey. *arXiv preprint arXiv:2411.18279*, 2024a.
- Jiwen Zhang, Jihao Wu, Yihua Teng, Minghui Liao, Nuo Xu, Xiao Xiao, Zhongyu Wei, and Duyu Tang. Android in the zoo: Chain-of-action-thought for gui agents, 2024b. <https://arxiv.org/abs/2403.02713>.
- Miaosen Zhang, Ziqiang Xu, Jialiang Zhu, Qi Dai, Kai Qiu, Yifan Yang, Chong Luo, Tianyi Chen, Justin Wagle, Tim Franklin, and Baining Guo. Phi-ground tech report: Advancing perception in gui grounding, 2025a. <https://arxiv.org/abs/2507.23779>.
- Zhong Zhang, Yaxi Lu, Yikun Fu, Yupeng Huo, Shenzhi Yang, Yesai Wu, Han Si, Xin Cong, Haotian Chen, Yankai Lin, Jie Xie, Wei Zhou, Wang Xu, Yuanheng Zhang, Zhou Su, Zhongwu Zhai, Xiaoming Liu, Yudong Mei, Jianming Xu, Hongyan Tian, Chongyi Wang, Chi Chen, Yuan Yao, Zhiyuan Liu, and Maosong Sun. Agentcpm-gui: Building mobile-use agents with reinforcement fine-tuning, 2025b. <https://arxiv.org/abs/2506.01391>.
- Yaowei Zheng, Junting Lu, Shenzhi Wang, Zhangchi Feng, Dongdong Kuang, and Yuwen Xiong. Easyr1: An efficient, scalable, multi-modality rl training framework. <https://github.com/hiyouga/EasyR1>, 2025.
- Zhipu-AI. Glm-4.5v. Available at: <https://docs.z.ai/guides/vlm/glm-4.5v>, 2025.
- Yuqi Zhou, Sunhao Dai, Shuai Wang, Kaiwen Zhou, Qinglin Jia, and Jun Xu. Gui-g1: Understanding r1-zero-like training for visual grounding in gui agents. 2025. <https://arxiv.org/abs/2505.15810>.
- Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, Zhangwei Gao, Erfei Cui, Xuehui Wang, Yue Cao, Yangzhou Liu, Xingguang Wei, Hongjie Zhang, Haomin Wang, Weiyue Xu, Hao Li, Jiahao Wang, Nianchen Deng, Songze Li, Yinan He, Tan Jiang, Jiapeng Luo, Yi Wang, Conghui He, Botian Shi, Xingcheng Zhang, Wenqi Shao, Junjun He, Yingtong Xiong, Wenwen Qu, Peng Sun, Penglong Jiao, Han Lv, Lijun Wu, Kaipeng Zhang,

Huipeng Deng, Jiaye Ge, Kai Chen, Limin Wang, Min Dou, Lewei Lu, Xizhou Zhu, Tong Lu, Dahua Lin, Yu Qiao, Jifeng Dai, and Wenhui Wang. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models, 2025.  
<https://arxiv.org/abs/2504.10479>.

## A Prompt Templates

### A.1 Grounding

#### Grounding Prompt

Outline the position corresponding to the instruction: {problem}. The output should be only [x1,y1,x2,y2].

### A.2 Navigation

#### Navigation Prompt

\*\*You are a GUI Agent\*\*.

Your task is to analyze a given user task, review current screenshot and previous actions, and determine the next action to complete the task.

### User Task

{problem}

### Previous Actions

{history}

### Available Actions

You may execute one of the following functions:

Click(box=(x1, y1))

Drag(start=(x1, y1), end=(x2, y2))

Scroll(start=(x1, y1), end=(x2, y2), direction='down/up/right/left')

Type(content=“”)

Launch(app=“”)

Wait()

Finished(content=“”)

CallUser(content=“”)

LongPress(box=(x1, y1))

PressBack()

PressHome()

PressEnter()

PressRecent()

### Instruction

- Make sure you understand the task goal to avoid wrong actions.
- Make sure you carefully examine the current screenshot. Sometimes the summarized history might not be reliable, over-claiming some effects.
- For requests that are questions (or chat messages), remember to use the ‘CallUser’ action to reply to user explicitly before finishing! Then, after you have replied, use the Finished action if the goal is achieved.
- Consider exploring the screen by using the ‘scroll’ action with different directions to reveal additional content.
- To copy some text: first select the exact text you want to copy, which usually also brings up

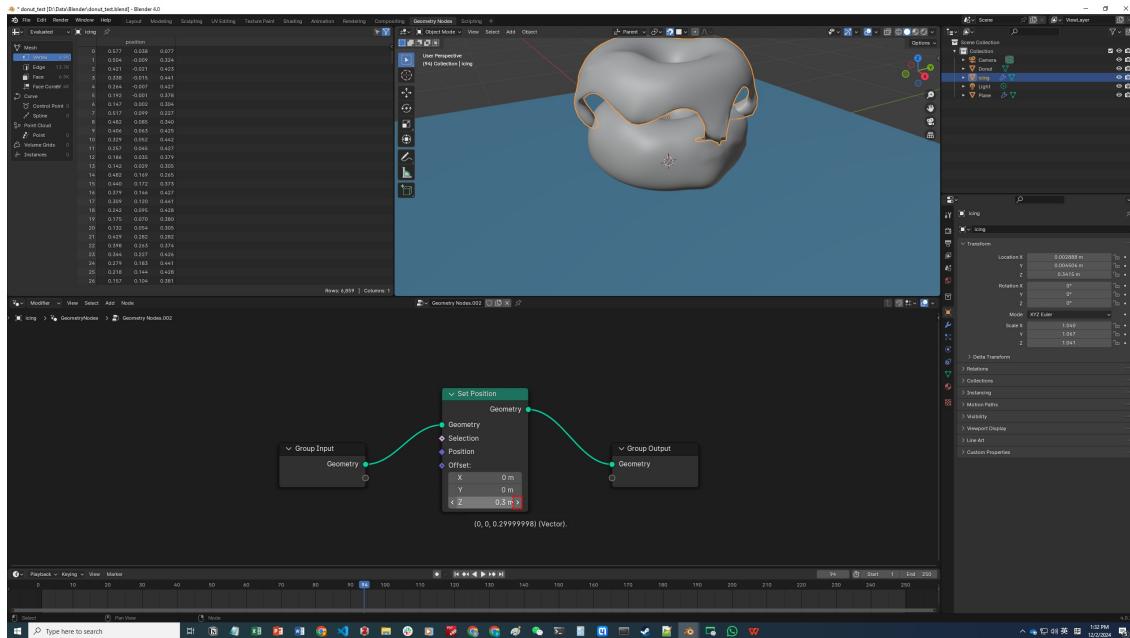
the text selection bar, then click the ‘copy’ button in bar.

- To paste text into a text box, first long press the text box, then usually the text selection bar will appear with a ‘paste’ button in it.
- You first think about the reasoning process in the mind, then provide the action. The reasoning and action are enclosed in <think></think> and <action></action> tags respectively. After providing action, summarize your action in <conclusion></conclusion> tags.

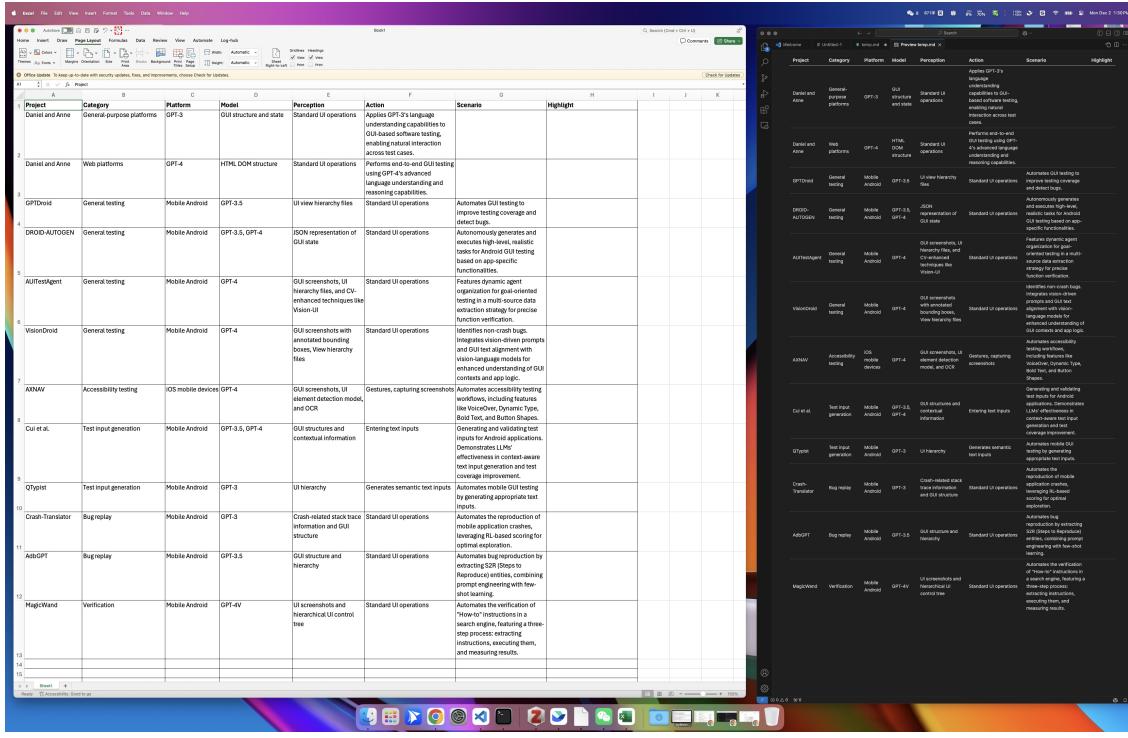
## B Qualitative Examples

### B.1 Grounding

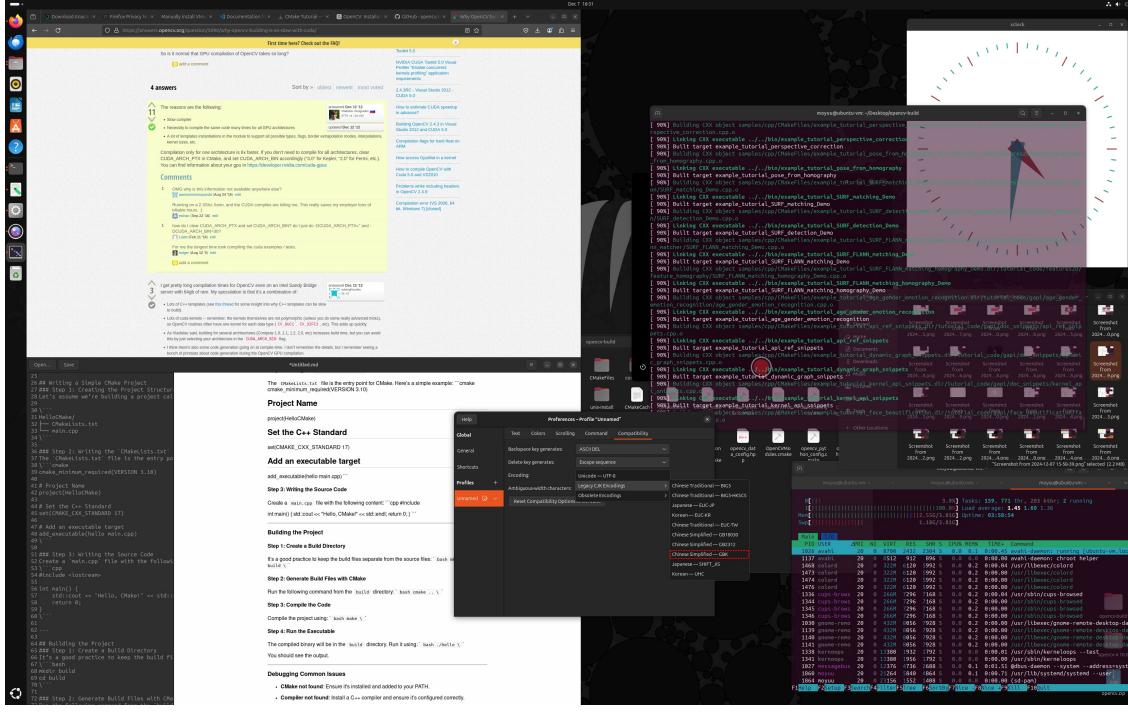
We show several grounding examples across desktop professional software, desktop system, website and mobile in Figure 5-10.



**Figure 5** Grounding example for desktop professional software Blender. The instruction is "Increase Z axis" and the bounding box result is displayed as a dotted red box



**Figure 6** Grounding example for Excel. The instruction is "Redo" and the bounding box result is displayed as a dotted red box.



**Figure 7** Grounding example for common linux system. The instruction is "Change the terminal encoding to the legacy GBK" and the bounding box result is displayed as a dotted red box.

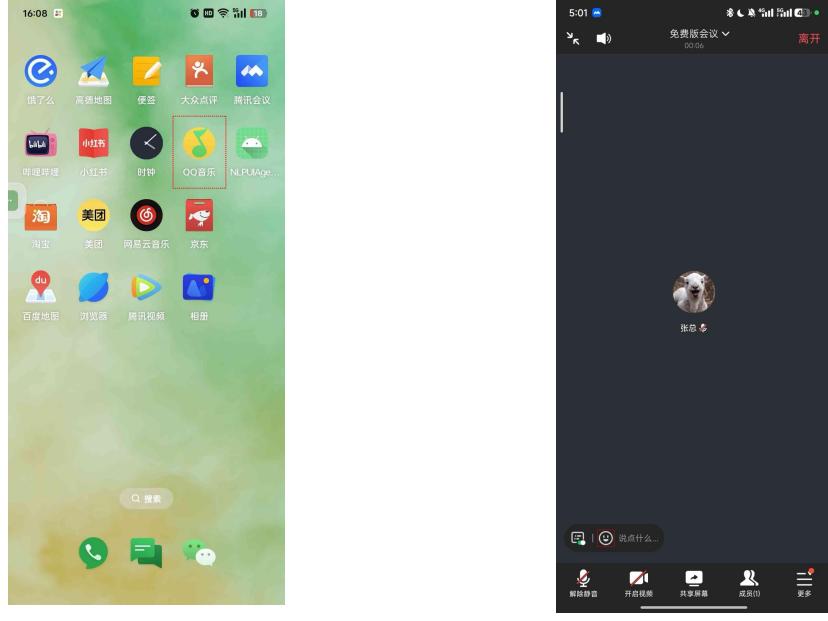
The screenshot shows the GitLab 'To-Do List' page. At the top, there's a navigation bar with a search bar and various icons. Below it, the title 'To-Do List' is displayed. A header bar includes filters for 'Group', 'Project', 'Author', 'Type', 'Action', and 'Last created'. The main area lists several tasks:

- Autocomplete interface guideline updates (#267) - Primer / design !455 (1 day ago, Could not merge)
- 404 for many URLs - x-lab / a11yproject.com #1 (Due Jan 3, 2030 - 1 day ago, You assigned to yourself.)
- Clarify usage of flash alert - Primer / design #316 (2 days ago, You assigned to yourself.)
- Feature/replace gulp - The A11Y Project / a11yproject.com !1537 (2 days ago, You requested a review from yourself.)
- Draft: Resolve "Feature: Search bar" - The A11Y Project / a11yproject.com !1535 (2 days ago, You requested a review from yourself.)
- Draft: Redesign - The A11Y Project / a11yproject.com !1534 (2 days ago, You requested a review from yourself.)
- Debug build time - Primer / design !454 (3 days ago, You requested a review from yourself.)

**Figure 8** Grounding example for the website Gitlab. The instruction is "Add a new one" and the bounding box result is displayed as a dotted red box.

The screenshot shows the Magento Business Intelligence service interface. The left sidebar has a navigation menu with icons and labels: DASHBOARD, SALES, CATALOG (highlighted with a red dotted box), CUSTOMERS (highlighted with a red dotted box), MARKETING, CONTENT, REPORTS, STORES, SYSTEM, and FIND PARTNERS & EXTENSIONS. The main content area is titled 'Catalog' and displays a message about synchronizing data with the Magento BI service. It includes a 'System Messages: 1' notification, a search bar, and a user dropdown for 'admin'. A large button labeled 'Reload Data' is visible. Below this, there's a section for reporting business performance with tabs for 'Orders' and 'Amounts', a date selector set to 'Today', and a message 'No Data Found'. At the bottom, there are summary statistics: Revenue \$0.00, Tax \$0.00, Shipping \$0.00, and Quantity 0. A footer navigation bar includes links for Bestsellers, Most Viewed Products, New Customers, and Customers.

**Figure 9** Grounding example for a shopping website. The instruction is "Switch to customers" and the bounding box result is displayed as a dotted red box.



Instruction: "QQ音乐"

Instruction: "该笑脸图标代表一个聊天界面中的表情"

**Figure 10** Grounding examples for the mobile system with Chinese instructions. The bounding box results are displayed as dotted red boxes.

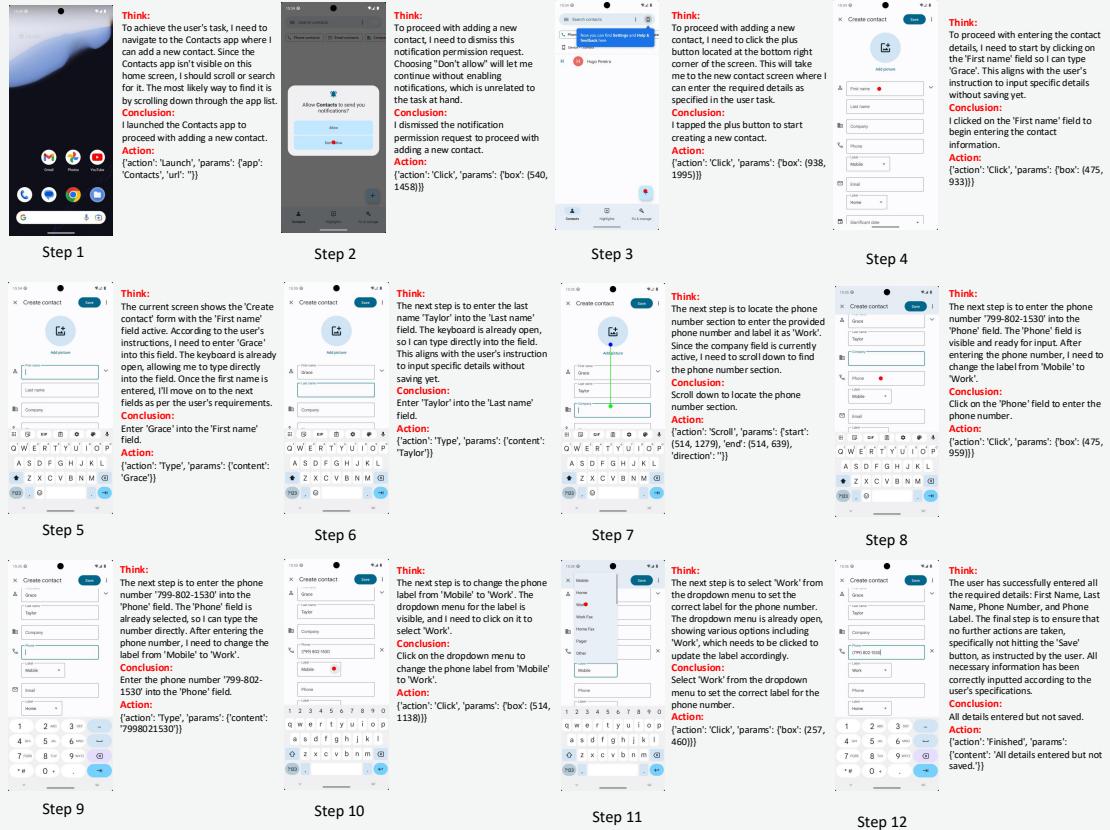
## B.2 Navigation

We present several UI navigation case studies in Figures 12 and 13 from the AndroidWorld dataset. To demonstrate UI-Venus’s cross-lingual capabilities, we additionally provide a navigation trace (Figure 11) showcasing its performance on a Chinese language application with Chinese interface elements.



**Figure 11** One UI-Venus navigation trace on a Chinese application with Chinese goal description. Interaction points are highlighted using red circular markers.

**Goal:** Go to the new contact screen and enter the following details: First Name: Grace, Last Name: Taylor, Phone: 799-802-1530, Phone Label: Work. Do NOT hit save.



**Figure 12** One trace of UI-Venus on the task named *ContactsNewContactDraft* in AndroidWorld. We mark the clicking points with red circles and observe that UI-Venus successfully create the draft without hitting the Save, which shows the model’s powerful navigation generalization and strong instruction-following ability.



**Figure 13** One trace of UI-Venus on the task named *MarkorDeleteAllNotes* in AndroidWorld. We can observe that UI-Venus successfully achieves the goal and has the reflection ability in Step 3. However, there also exists the conflict between think and action in Step 5, remaining as a future work about how to solve MLLM's hallucination.