# UI-R1: Enhancing Action Prediction of GUI Agents by Reinforcement Learning

**Zhengxi Lu**[1†], **Yuxiang Chai**[2†], **Yaxuan Guo**[1], **Xi Yin**[1],
**Liang Liu**[1‡], **Hao Wang**[1], **Guanjing Xiong**[1], **Hongsheng Li**[2✉]

[1] vivo AI Lab    [2] MMLab @ CUHK
[†] Equal Contribution    [‡] Project Lead    [✉] Corresponding Author

## Abstract

The recent DeepSeek-R1 has showcased the emergence of reasoning capabilities in LLMs through reinforcement learning (RL) with rule-based rewards. Building on this idea, we are the first to explore how rule-based RL can enhance the reasoning capabilities of multimodal large language models (MLLMs) for graphic user interface (GUI) action prediction tasks. To this end, we curate a small yet high-quality dataset of 136 challenging tasks, encompassing five common action types on mobile devices. We also introduce a unified rule-based action reward, enabling model optimization via policy-based algorithms such as Group Relative Policy Optimization (GRPO). Experimental results demonstrate that our proposed data-efficient model, **UI-R1-3B**, achieves substantial improvements on both in-domain (ID) and out-of-domain (OOD) tasks. Specifically, on the ID benchmark ANDROIDCONTROL, the action type accuracy improves by **15%**, while grounding accuracy increases by **10.3%**, compared with the base model (i.e. Qwen2.5-VL-3B). On the OOD GUI grounding benchmark ScreenSpot-Pro, our model surpasses the base model by **6.0%** and achieves competitive performance with larger models (e.g., OS-Atlas-7B), which are trained via supervised fine-tuning (SFT) on 76K data. These results underscore the potential of rule-based reinforcement learning to advance GUI understanding and control, paving the way for future research in this domain.
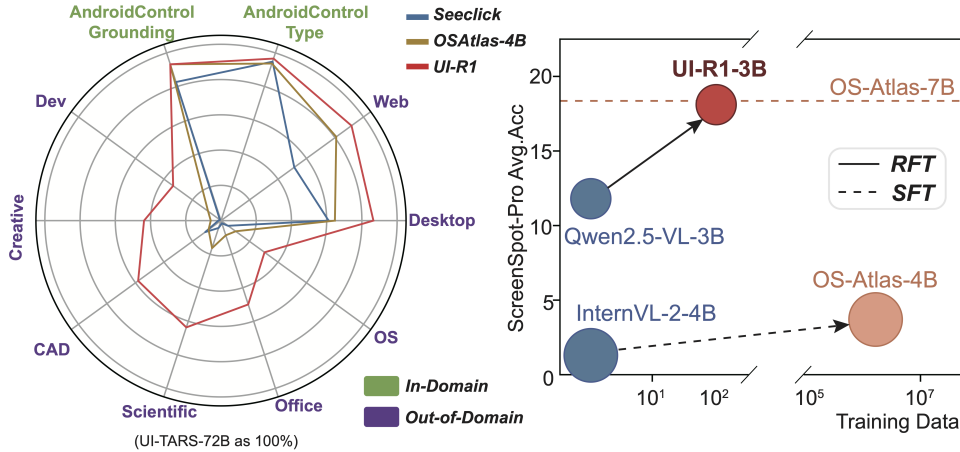
Figure 1: **Left**: Overall performance of UI-R1-3B on both in-domain (i.e., ANDROIDCONTROL) and out-of-domain (i.e., ScreenSpot-Pro, ScreenSpot desktop and web subsets) tasks; **Right**: Employing reinforcement fine-tuning (RFT), UI-R1-3B achieves performance comparable to SFT models with significantly fewer data and GPU hours. The circle radius indicates the model size.

# 1   Introduction

Supervised fine-tuning (SFT) has long been the standard training paradigm for large language models (LLMs) and graphic user interface (GUI) agents (Qin et al., 2025; Wu et al., 2024; Hong et al., 2024). However, SFT relies heavily on large-scale, high-quality labeled datasets, leading to prolonged training times and high computational costs. Furthermore, existing open-source VLM-based GUI agents trained using SFT can be criticized for poor performance in out-of-domain (OOD) scenarios (Lu et al., 2024; Chai et al., 2024), limiting their effectiveness and applicability in real-world applications.

Rule-based reinforcement learning (RL) or reinforcement fine-tuning (RFT) has recently emerged as an efficient and scalable alternative to SFT for the development of LLMs, which efficiently fine-tune the model with merely dozens to thousands of samples to excel at domain-specific tasks. It uses predefined task-specific reward functions, eliminating the need for costly human annotations. Recent works, such as DeepSeek-R1 (Guo et al., 2025), demonstrate the effectiveness of rule-based RL in mathematical problem solving by evaluating the correctness of the solution, while others (Liu et al., 2025; Wang et al., 2025; Peng et al., 2025; Chen et al., 2025) extend the algorithm to multimodal models, achieving notable improvements in vision-related tasks such as image grounding and object detection. By focusing on measurable objectives, rule-based RL enables practical and versatile model optimization across both textual and multimodal domains, offering significant advantages in terms of efficiency, scalability, and reduced reliance on large datasets.

Previous studies targeting traditional vision-related tasks always rely on the traditional Intersection over Union (IoU) metric commonly used for grounding and detection tasks. In this work, we extend the rule-based RL paradigm to a new application domain by focusing on GUI action prediction tasks driven by low-level instructions. To achieve this, MLLM generates multiple responses (trajectories) that contain the reasoning tokens and the final answers for each input. Then our proposed reward function evaluates each response and updates the model by policy optimization, such as GRPO (Shao et al., 2024), to improve its reasoning ability. Our reward function contains the action type reward, the action argument reward, along with the common format reward. In detail, (1) the action type reward is determined by whether the predicted action type matches the ground truth; (2) the action argument reward (focused on `Click`), is evaluated by whether the predicted coordinates fall within the ground truth bounding box; (3) the format reward is evaluated by whether the model provides both the reasoning process and the final answer. This flexible and effective reward mechanism is well aligned with the objectives of general GUI-related tasks, ensuring both accuracy and interpretability in the model's performance.

Regarding data preparation, we follow Muennighoff et al. (2025) and select just 130+ training mobile samples according to three criterion: difficulty, diversity, and quality, making our method remarkably data-efficient. Experiments demonstrate that **UI-R1** achieves significant performance improvements on out-of-domain (OOD) data, including entries from desktop and web platforms, indicating the potential of rule-based RL to tackle complex GUI-related tasks across diverse domains effectively.

In summary, our contributions are as follows.

- We propose **UI-R1**, which enhances MLLM's reasoning capabilities on GUI action prediction tasks through DeepSeek R1 style reinforcement learning.

- We design a unified rule-based action reward function that effectively aligns with the objectives of common GUI tasks.

- We utilize the three-stage data selection method and collect only 130+ high-quality training data from the mobile domain. Despite limited data, our model **UI-R1-3B** achieves notable performance gains on out-of-domain benchmarks, such as those from desktop and web platforms, showcasing adaptability and generalization capability in GUI-related tasks.
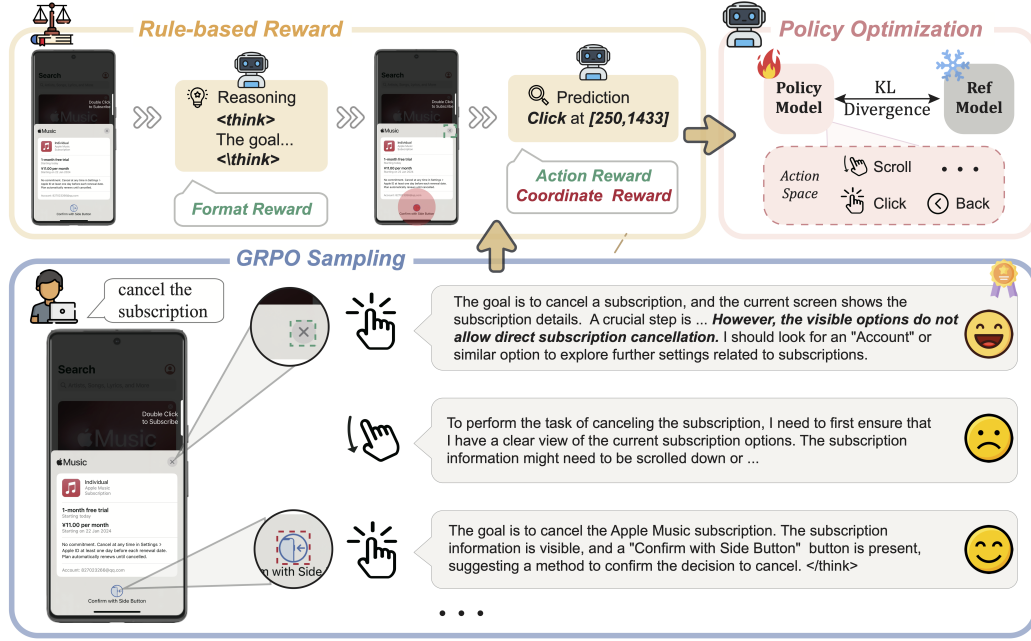
Figure 2: Overview of UI-R1 training framework. Given a GUI screenshot and a text instruction from the user, the policy model (i.e., Qwen2.5-VL-3B) generates multiple action planning responses with reasoning. Our proposed rule-based action reward function is then applied, and the policy model is updated using a policy gradient optimization algorithm.

## 2 Related Work

### 2.1 GUI Agents

Starting with CogAgent (Hong et al., 2024), researchers have used MLLMs for GUI-related tasks, including device control, task completion, GUI understanding, and more. One line of work, such as the AppAgent series (Zhang et al., 2023; Li et al., 2024b) and the Mobile-Agent series (Wang et al., 2024b;a), integrates commercial generalist models like GPT for planning and prediction tasks. These agents rely heavily on prompt engineering and multi-agent collaboration to execute complex tasks, making them adaptable but dependent on careful manual design for optimal performance. Another branch of research focuses on fine-tuning smaller open-source MLLMs on task-specific GUI datasets (Rawles et al., 2023; Li et al., 2024a; Chai et al., 2024; Gou et al., 2024) to create specialist agents. For example, Chai et al. (2024) enhances agents by incorporating additional functionalities of the GUI element in the Android system, while UGround(Gou et al., 2024) develops a special GUI grounding model tailored for precise GUI element localization. Wu et al. (2024) develops a foundational model for GUI action prediction. Moving beyond task-specific fine-tuning, UI-TARs (Qin et al., 2025) introduces a more comprehensive approach by combining GUI-related pretraining with task-wise reasoning fine-tuning, aiming to better align models with the intricacies of GUI interactions. Despite their differences, all of these existing agents share a common reliance on the SFT paradigm. This training approach, while effective, depends heavily on large-scale, high-quality labeled datasets.

### 2.2 Rule-Based Reinforcement Learning

Rule-based reinforcement learning (RL) has recently emerged as an efficient alternative to traditional training paradigms by leveraging predefined rule-based reward functions to guide model behavior. DeepSeek-R1 (Guo et al., 2025) first introduced this approach, using

reward functions based on predefined criteria, such as checking whether an LLM's final answer matches the ground truth for math problems. The reward focuses solely on the final results, leaving the reasoning process to be learned by the model itself. Zeng et al. (2025) reproduces the algorithm on models with smaller sizes and illustrates its effectiveness on small language models. Subsequent works (Chen et al., 2025; Shen et al., 2025; Liu et al., 2025; Wang et al., 2025; Peng et al., 2025; Meng et al., 2025), extended the paradigm to multimodal models by designing task-specific rewards for visual tasks, including correct class predictions for image classification and IoU metrics for image grounding and detection. These studies demonstrate the adaptability of rule-based RL for both pure-language and multimodal models. By focusing on task-specific objectives without requiring extensive labeled datasets or human feedback, rule-based RL shows strong potential as a scalable and effective training paradigm across diverse tasks.

## 3 Method

UI-R1 is a reinforcement learning (RL) training paradigm designed to enhance a GUI agent's ability to successfully complete low-level instructional tasks. We define "low-level instructions" as directives that guide the agent to perform actions based on a single state (e.g., a GUI screenshot), consistent with the definition in ANDROIDCONTROL (Li et al., 2024a). For example, *"Click the menu icon in the top left corner"* represents a low-level instruction, whereas *"Create an event for 2 PM tomorrow"* is a high-level instruction. The specifics of the training data selection and reward function design are detailed in the following sections. Figure 2 illustrates the main parts of the framework.

### 3.1 Preliminary

Many rule-based RL works (Guo et al., 2025; Zeng et al., 2025; Liu et al., 2025) adopt the Group Relative Policy Optimization (GRPO) algorithm (Shao et al., 2024) for RL training. GRPO offers an alternative to commonly used Proximal Policy Optimization (PPO) (Schulman et al., 2017) by eliminating the need for a critic model. Instead, GRPO directly compares a group of candidate responses to determine their relative quality.

In GRPO, given a task question, the model generates a set of $N$ potential responses $\{o_1, o_2, \ldots, o_N\}$. Each response is evaluated by taking the corresponding actions and computing its reward $\{r_1, r_2, \ldots, r_N\}$. Unlike PPO, which relies on a single reward signal and a critic to estimate the value function, GRPO normalizes these rewards to calculate the relative advantage of each response. The relative quality $A_i$ of the $i$-th response is computed as

$$A_i = \frac{r_i - Mean(\{r_1, r_2, \ldots, r_N\})}{Std(\{r_1, r_2, \ldots, r_N\})}, \tag{1}$$

where *Mean* and *Std* represent the mean and standard deviation of the rewards, respectively. This normalization step ensures that responses are compared within the context of the group, allowing GRPO to better capture nuanced differences between candidates. Policy updates are further constrained by minimizing the KL divergence between the updated and reference models, ensuring stable RL learning.

### 3.2 Rule-Based Action Rewards

The rule-based reward function introduced by DeepSeek-R1 (Guo et al., 2025) represents a foundational step in rule-based RL by simply evaluating whether model predictions exactly match ground-truth answers. This straightforward approach efficiently aligns models with preference alignment algorithms and provides clear optimization signals. For vision-related tasks, works such as VLM-R1 (Shen et al., 2025) and Visual-RFT (Liu et al., 2025) extend this idea by designing task-specific rewards. For image grounding tasks, they compute the IoU between the predicted and ground-truth bounding boxes as the reward. Similarly, for image classification tasks, rewards are determined by checking whether the predicted and ground-truth classes match.

In GUI-related tasks, the ability of GUI grounding and understanding is a critical requirement for agents. Unlike traditional image grounding tasks, GUI grounding requires agents to identify where specific actions, such as `click`, should be performed on a given GUI screenshot. To address this unique gap, we propose a reward function tailored for GUI tasks, as defined in Equation 2:

$$R_{\mathcal{A}} = R_{\mathcal{T}} + R_{\mathcal{C}} + R_{\mathcal{F}}, \tag{2}$$

where the predicted action $\mathcal{A} = \{\mathcal{T}, \mathcal{C}\}$ consists of two components: $\mathcal{T}$, which represents the action type (e.g., `click`, `swipe`), and $\mathcal{C}$, which represents the `click` coordinate. $R_{\mathcal{F}}$ represents the common response format reward.

**Action type reward.** In our tasks, the action space includes `Click`, `Scroll`, `Back`, `Open_App`, and `Input_Text`, covering a wide range of common application scenarios in daily life, as inspired by GUIPivot (Wu et al., 2025). The action type reward, denoted as $R_{\mathcal{T}}$, is computed by comparing the predicted action type $\mathcal{T}'$ with the ground truth action type $\mathcal{T}$. It assigns a reward of 1 if $\mathcal{T}' = \mathcal{T}$ and 0 otherwise, providing a straightforward and effective evaluation mechanism for action type prediction.

**Coordinate accuracy reward.** Through observation, we find that among all action types, the most common action argument error occurs in the mis-prediction of coordinates for the `click` action when given a low-level instruction. To address this issue, we specifically design a coordinate accuracy reward. The model is required to output a coordinate $\mathcal{C} = [x, y]$, indicating where the `click` action should be performed. Given the ground truth bounding box $\mathcal{B} = [x1, y1, x2, y2]$, the coordinate accuracy reward $R_{\mathcal{C}}$ is computed as shown in Equation 3:

$$R_{\mathcal{C}} = \begin{cases} 1 & \text{if coord } \mathcal{C} \text{ } in \text{ box } \mathcal{B}, \\ 0 & \text{else.} \end{cases} \tag{3}$$

Unlike general visual grounding tasks which compute the IoU between the predicted bounding box and the ground truth box, our approach prioritizes action coordinate prediction over element grounding. This focus is more appropriate for GUI agents and better aligns with human intuition, as the ultimate goal is to ensure correct actions are performed rather than merely locating GUI elements.

**Format reward.** During training, we incorporate the widely-used format reward to guide the model in generating its reasoning process and final answer in a structured format. This decision is based on our simple experiment that agents producing reasoning processes outperform those directly outputting action predictions by approximately 6%. The reasoning process plays a key role in the model's self-learning and iterative improvement during reinforcement fine-tuning, while the reward tied to the final answer drives optimization. The format reward, denoted as $R_{\mathcal{F}}$, ensures that the model's predictions follow the required HTML tag format, specifically using <*think*> for the reasoning process and <*answer*> for the final answer. This structured output not only enhances clarity, but also ensures consistency in the model's predictions.

> **Prompt for Training and Inference**
>
> In this GUI screenshot, I want to perform the command ***instruction***. Please provide the action to perform (enumerate in [click, open_app, scroll, navigate_back, input_text]) and the coordinate where the cursor is moved to(integer) if click is performed. Output the thinking process in <*think*> </*think*> and final answer in <*answer*> </*answer*> tags. The output answer format should be as follows: <*think*> ... </*think*> <*answer*>[action: enum[click, open_app, scroll, navigate_back, input_text], coordinate: [x, y]]</*answer*>. Please strictly follow the format.

### 3.3 Training Data Selection

Compared to SFT, rule-based RL has demonstrated the capability to achieve comparable or even superior performance on mathematical and vision-related tasks using only a limited number of training samples (Zeng et al., 2025; Liu et al., 2025). Building on this efficiency and inspired by s1 (Muennighoff et al., 2025), we implement a three-stage data selection process to refine open-source GUI-related datasets based on three key principles: Quality, Difficulty, and Diversity. The detailed distribution of the dataset can be found in Appendix A.2.

**Quality.** For refining the `click` action arguments, we use the mobile subset of ScreenSpot (Cheng et al., 2024) as our initial dataset. ScreenSpot offers clean and well-aligned task-element paired annotations, making it ideal for defining and calculating $R_C$. For other actions, we randomly select 1K episodes from ANDROIDCONTROL (Li et al., 2024a), as it shares a similar action space and provides low-level instructions. However, since the element annotations in ANDROIDCONTROL are unfiltered and misaligned, we exclude `click` action steps and retain the rest.

**Difficulty.** To identify hard samples, we evaluated Qwen2.5-VL-3B on each task instruction by model performance, where a sample is labeled "hard" if the model's output does not match the ground truth. We only keep the "hard" samples among all the data collected.

**Diversity.** We ensure diversity by selecting samples with different action types in ANDROIDCONTROL (e.g., `Scroll`, `Back`, `Open App`, `Input Text`) and element types in ScreenSpot (e.g. Icon, Text). Rare actions, such as `Wait` and `Long Press`, are excluded from ANDROIDCONTROL. After applying these criteria, we finalize a high-quality mobile training dataset consisting of 136 samples.

## 4 Experiment

### 4.1 GUI Grounding Capability

We assess the grounding capability of UI-R1 using two benchmarks: ScreenSpot (Cheng et al., 2024) and ScreenSpot-Pro (Li et al., 2025). ScreenSpot evaluates GUI grounding capability across mobile, desktop, and web platforms, while ScreenSpot-Pro focuses on high-resolution professional environments, featuring expert-annotated tasks spanning 23 applications, five industries, and three operating systems. Evaluation results of ScreenSpotV2 (Wu et al., 2024) are in Appendix B.

| Model | Method | Model Size | Data Size | Web Icon | Web Text | Desktop Icon | Desktop Text | Average |
|---|---|---|---|---|---|---|---|---|
| **Supervised Fine-tuning** | | | | | | | | |
| SeeClick | SFT | 9.6B | 1M | 32.5 | 55.7 | 30.0 | 72.2 | 49.0 |
| CogAgent | SFT | 18B | - | 28.6 | 70.4 | 20.0 | 74.2 | 51.0 |
| Qwen2.5-VL | SFT | 3B | 500 | 63.1 | 78.3 | 46.4 | 85.0 | 70.1 |
| UGround-V1 | SFT | 7B | 10M | 70.4 | 80.4 | <u>63.6</u> | 82.5 | 75.2 |
| AGUVIS | SFT | 7B | 1M | <u>70.7</u> | **88.1** | **74.8** | <u>85.7</u> | **80.4** |
| **Zero Shot / Reinforcement Learning** | | | | | | | | |
| Qwen2-VL | ZS | 7B | 0 | 25.7 | 35.2 | 54.3 | 76.3 | 46.5 |
| Qwen2.5-VL | ZS | 3B | 0 | 43.2 | 60.0 | 40.0 | 80.9 | 57.1 |
| UI-R1 | RFT | **3B** | 136 | **73.3** | <u>85.2</u> | 59.3 | **90.2** | <u>78.6</u> |

Table 1: Grounding accuracy on ScreenSpot. The optimal and the suboptimal results are **bolded** and <u>underlined</u>, respectively. ZS indicates zero-shot OOD inference and RFT indicates rule-based reinforcement learning.

| Model | Development | | Creative | | CAD | | Scientific | | Office | | OS | | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Text | Icon | Text | Icon | Text | Icon | Text | Icon | Text | Icon | Text | Icon | |
| **Supervised Fine-tuning** | | | | | | | | | | | | | |
| SeeClick | 0.6 | 0.0 | 1.0 | 0.0 | 2.5 | 0.0 | 3.5 | 0.0 | 1.1 | 0.0 | 2.8 | 0.0 | 1.1 |
| OS-Atlas-4B | 7.1 | 0.0 | 3.0 | 1.4 | 2.0 | 0.0 | 9.0 | 5.5 | 5.1 | 3.8 | 5.6 | 0.0 | 3.7 |
| ShowUI-2B | 16.9 | 1.4 | 9.1 | 0.0 | 2.5 | 0.0 | 13.2 | 7.3 | 15.3 | 7.5 | 10.3 | 2.2 | 7.7 |
| CogAgent-18B | 14.9 | 0.7 | 9.6 | 0.0 | 7.1 | 3.1 | 22.2 | 1.8 | 13.0 | 0.0 | 5.6 | 0.0 | 7.7 |
| Aria-GUI | 16.2 | 0.0 | 23.7 | 2.1 | 7.6 | 1.6 | 27.1 | 6.4 | 20.3 | 1.9 | 4.7 | 0.0 | 11.3 |
| Qwen2.5-VL-3B* | 15.6 | 0.7 | 13.1 | 2.1 | 5.6 | 3.1 | 27.8 | 8.1 | 20.3 | 5.7 | 14.0 | 0.0 | 10.8 |
| UGround-7B | 26.6 | 2.1 | 27.3 | 2.8 | <u>14.2</u> | 1.6 | 31.9 | 2.7 | 31.6 | 11.3 | <u>17.8</u> | 0.0 | 16.5 |
| Claude** | 22.0 | <u>3.9</u> | 25.9 | <u>3.4</u> | **14.5** | 3.7 | 33.9 | **15.8** | 30.1 | 16.3 | 11.0 | **4.5** | 17.1 |
| OS-Atlas-7B | **33.1** | 1.4 | 28.8 | 2.8 | 12.2 | <u>4.7</u> | <u>37.5</u> | 7.3 | 33.9 | 5.7 | **27.1** | 4.5 | **18.9** |
| **Zero Shot / Reinforcement Fine-tuning** | | | | | | | | | | | | | |
| Qwen-VL-7B | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.7 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.1 |
| GPT-4o | 1.3 | 0.0 | 1.0 | 0.0 | 2.0 | 0.0 | 2.1 | 0.0 | 1.1 | 0.0 | 0.0 | 0.0 | 0.8 |
| Qwen2-VL-7B | 2.6 | 0.0 | 1.5 | 0.0 | 0.5 | 0.0 | 6.3 | 0.0 | 3.4 | 1.9 | 0.9 | 0.0 | 1.6 |
| Qwen2.5-VL-3B | 14.9 | 2.1 | 20.2 | 1.4 | 4.1 | <u>4.7</u> | 34.0 | 7.3 | 22.0 | 3.8 | 6.5 | 2.2 | 11.8 |
| UI-R1-3B | 22.7 | **4.1** | <u>27.3</u> | **3.5** | 11.2 | **6.3** | **42.4** | <u>11.8</u> | <u>32.2</u> | <u>11.3</u> | 13.1 | **4.5** | <u>17.8</u> |

Table 2: Accuracy on ScreenSpot-Pro. The optimal and the suboptimal results are **bolded** and <u>underlined</u>, respectively. * Qwen2.5-VL-3B here is supervised fine-tuned on 500 ScreenSpot-mobile data. ** Claude refers to *Claude-computer-use*.

**Setting**  We train the Qwen2.5-VL-3B model on the three-stage selected data (details in Section 3.3) using rule-based RL, naming the resulting model UI-R1-3B. Furthermore, we train the base model using supervised fine-tuning on the entire ScreenSpot mobile set, referring to it as Qwen2.5-VL-3B* in Table 2. For evaluation, an action prediction is considered correct if the predicted `click` coordinate lies within the ground truth bounding box. Accuracy is computed as the ratio of the correct predictions to the total number of test samples.

**Analysis**  Experimental results show that our method significantly improves the GUI grounding capability of the 3B model (**+20%** on ScreenSpot and **+6%** on ScreenSpot-Pro from Table 1 and Table 2), surpassing most 7B models on both benchmarks. Additionally, it also achieves performance comparable to the SOTA 7B models (i.e. AGUVIS (Xu et al., 2024) and OS-Atlas (Wu et al., 2024)), which are trained using supervised fine-tuning on substantially larger labeled grounding datasets.

Qwen2.5-VL-3B (SFT) in Table 1 demonstrates that supervised fine-tuning (SFT) with a limited amount of data (e.g., 500 samples) can effectively improve in-domain performance by tailoring the model to specific tasks. However, the comparison between Qwen2.5-VL-3B (ZS) and Qwen2.5-VL-3B (SFT) in Table 2 highlights a critical limitation of SFT: its effectiveness significantly diminishes in OOD scenarios. This limitation arises from the dependency of SFT on task-specific labeled data, restricting the model's ability to adapt to unseen environments. In contrast, our RL approach not only enhances OOD generalization by focusing on task-specific reward optimization, but also achieves with far fewer training samples, offering a scalable and efficient alternative to traditional SFT methods.

## 4.2 Action Prediction Capability

We further evaluate the model's ability to predict single-step actions based on low-level instructions. As described in Section 3.3, we test our model on a selected subset of AN-DROIDCONTROL. The low-level instructions in ANDROIDCONTROL enrich the ScreenSpot benchmark by introducing a wider range of action types.

**Setting**  The accuracy of the action prediction is evaluated by the accuracies of action type and grounding: (1) The action type accuracy evaluates the match rate between the predicted action types (e.g., `click`, `scroll`) and ground truth types; (2) The grounding accuracy focuses specifically on the accuracy of `click` action argument predictions, similar to Section 4.1. Since ground truth bounding boxes are not consistently available in the

| Model | Method | Model size | Data size | Type | Grounding | Average |
|---|---|---|---|---|---|---|
| **Supervised Fine-tuning** | | | | | | |
| SeeClick | SFT | 9.6B | 76K | 93.0 | 73.4 | 83.2 |
| InternVL-2 | SFT | 4B | 76K | 90.9 | <u>84.1</u> | 87.5 |
| GUIPivot-Qwen | SFT | 7B | 76K | **96.8** | 75.1 | 86.0 |
| OS-Atlas | SFT | 4B | 76K | 91.9 | 83.8 | 87.8 |
| OS-Atlas | SFT | 7B | 76K | 93.6 | **88.0** | **90.8** |
| **Zero Shot / Reinforcement Fine-tuning** | | | | | | |
| GPT-4o | ZS | – | 0 | 74.3 | 38.7 | 56.5 |
| OS-Atlas | ZS | 4B | 0 | 64.6 | 71.2 | 67.9 |
| OS-Atlas | ZS | 7B | 0 | 73.0 | 73.4 | 73.2 |
| Qwen2.5-VL | ZS | 3B | 0 | 79.3 | 72.3 | 75.8 |
| UI-R1 | RFT | 3B | 136 | <u>94.3</u> | 82.6 | <u>88.5</u> |

Table 3: Low-level agent capabilities on ANDROIDCONTROL. The Average column computes the mean of Type and Grounding scores.

ANDROIDCONTROL test data, we measure performance by calculating the distance between the predicted and ground truth coordinates. A prediction is considered correct if it falls within 14% of the screen size from the ground truth, following the evaluation method of UI-TARS (Qin et al., 2025).

**Analysis** As shown in Table 3, the comparison between UI-R1 and the Qwen2.5-VL (ZS) model highlights the significant benefits of the RL training framework. UI-R1 improves the accuracy of action type prediction by **15%** and click element grounding accuracy by **20%**, all while using only 136 training data points. Compared with other SFT models, the evaluation results illustrate that UI-R1 not only excels in scenarios with extremely limited training data but also achieves superior type accuracy and grounding performance even than larger models. This underscores the effectiveness of the RL training framework in leveraging small datasets to achieve substantial performance gains, demonstrating its potential as a highly data-efficient and scalable approach for training models in resource-constrained environments.

### 4.3 Performance Factor Study

**Data Scale** In Figure 3 (left), we investigate the relationship between training data size and model performance and compare the two methods of selecting training data from the entire dataset: random selection and `select by difficulty` (as in Section 3.3). The second method involves selecting the top K tasks with the longest reasoning lengths that the Qwen2.5-VL-3B model fails to solve. We find that model performance improves as the training data size increases, but the trend is gradually saturating. Moreover, our `select by difficulty` method results in significantly better performance than random selection.

**Reasoning Length** In Figure 3 (right), the result reveals that as the reasoning length of the answer increases, the accuracy tends to decrease, suggesting that the questions are getting harder to answer. With reinforcement learning, UI-R1's reasoning ability is significantly enhanced, leading to more pronounced accuracy improvements, especially on more challenging questions.

### 4.4 Ablation Study

**Reward Function** The design of the reward function plays a crucial role in enabling the self-learning capabilities of the model. To evaluate this, we first examine the necessity of the two components of the reward function, `action + coord`. Specifically, the `action` reward improves action prediction accuracy, while the `coord` reward enhances the model's ability to ground `click` elements. Next, we compare this with an alternative reward design, `action + bbox`, where the coordinate reward $R_C$ is replaced by an IoU-based reward $R_{\text{IoU}}$
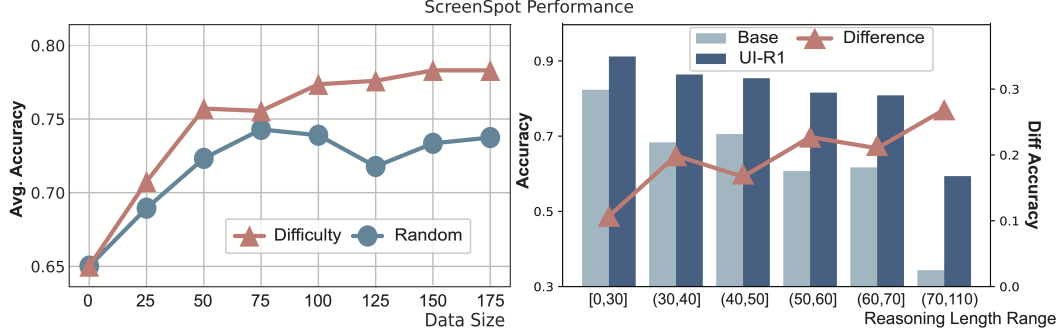
Figure 3: **Left**: Impact of data selection methods and data size; **Right**: Study of relation between answering accuracy and reasoning length.
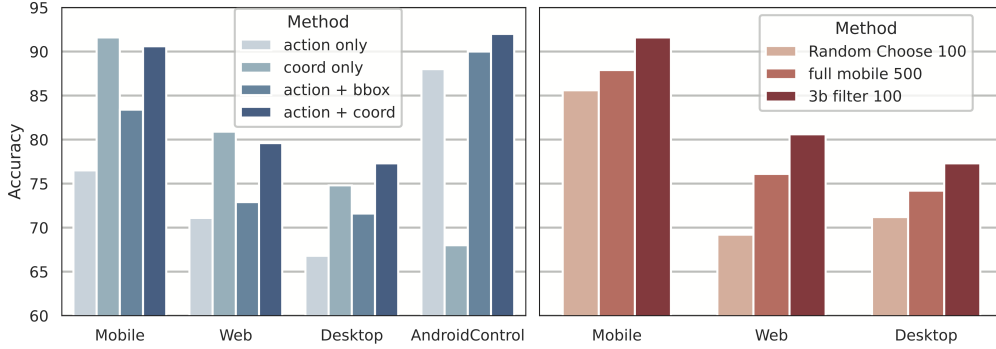


Figure 4: **Left**: Ablation on reward function; **Right**: Ablation on data selection method.

in Equation 2. In this setup, the IoU metric is calculated between the ground truth bounding box and the predicted box, and $R_{\text{IoU}}$ assigns a value of 1 if IoU $> 0.5$ and 0 otherwise.

Through ablation studies, as shown in Figure 4 (left), we demonstrate the superior effectiveness of $R_{\mathcal{C}}$ over $R_{\text{IoU}}$ for improving `click` element grounding. However, we also observe that the action reward does not always positively impact grounding tasks. This is likely because a larger action space can introduce ambiguity, making it harder for the model to focus solely on element grounding tasks. These findings highlight the importance of carefully balancing the reward components according to the specific objectives of the task.

**Data Selection**   We also examine the impact of different data selection methods, as shown in Figure 4 (right). A comparison of three methods across all domains demonstrates that neither random selection nor the use of the entire dataset matches the effectiveness of our three-stage data selection pipeline, indicating that the use of a smaller set of high-quality data can lead to higher performance.

## 5   Conclusion

We propose the UI-R1 framework, which extends rule-based reinforcement learning to GUI action prediction tasks, offering a scalable alternative to traditional Supervised Fine-Tuning (SFT). We designed a novel reward function that evaluates both action type and arguments, enabling efficient learning with reduced task complexity. Using only 130+ training samples from the mobile domain, our approach achieves significant performance improvements and strong generalization to out-of-domain datasets, including desktop and web platforms. The results demonstrate the adaptability, data efficiency, and ability of the rule-based RL to handle specialized tasks effectively.

9

# References

Yuxiang Chai, Siyuan Huang, Yazhe Niu, Han Xiao, Liang Liu, Dingyu Zhang, Peng Gao, Shuai Ren, and Hongsheng Li. Amex: Android multi-annotation expo dataset for mobile gui agents. *arXiv preprint arXiv:2407.17490*, 2024.

Liang Chen, Lei Li, Haozhe Zhao, Yifan Song, and Vinci. R1-v: Reinforcing super generalization ability in vision-language models with less than \$3. `https://github.com/Deep-Agent/R1-V`, 2025. Accessed: 2025-02-02.

Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Yantao Li, Jianbing Zhang, and Zhiyong Wu. Seeclick: Harnessing gui grounding for advanced visual gui agents. *arXiv preprint arXiv:2401.10935*, 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. *arXiv preprint arXiv:2410.05243*, 2024.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxiao Dong, Ming Ding, et al. Cogagent: A visual language model for gui agents. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14281–14290, 2024.

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, 2025.

Wei Li, William Bishop, Alice Li, Chris Rawles, Folawiyo Campbell-Ajala, Divya Tyamagundlu, and Oriana Riva. On the effects of data scale on computer control agents. *arXiv preprint arXiv:2406.03679*, 2024a.

Yanda Li, Chi Zhang, Wanqi Yang, Bin Fu, Pei Cheng, Xin Chen, Ling Chen, and Yunchao Wei. Appagent v2: Advanced agent for flexible mobile interactions. *arXiv preprint arXiv:2408.11824*, 2024b.

Ziyu Liu, Zeyi Sun, Yuhang Zang, Xiaoyi Dong, Yuhang Cao, Haodong Duan, Dahua Lin, and Jiaqi Wang. Visual-rft: Visual reinforcement fine-tuning. *arXiv preprint arXiv:2503.01785*, 2025.

Yadong Lu, Jianwei Yang, Yelong Shen, and Ahmed Awadallah. Omniparser for pure vision based gui agent. *arXiv preprint arXiv:2408.00203*, 2024.

Fanqing Meng, Lingxiao Du, Zongkai Liu, Zhixiang Zhou, Quanfeng Lu, Daocheng Fu, Botian Shi, Wenhai Wang, Junjun He, Kaipeng Zhang, et al. Mm-eureka: Exploring visual aha moment with rule-based large-scale reinforcement learning. *arXiv preprint arXiv:2503.07365*, 2025.

Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling, 2025. URL `https://arxiv.org/abs/2501.19393`.

Yingzhe Peng, Gongrui Zhang, Miaosen Zhang, Zhiyuan You, Jie Liu, Qipeng Zhu, Kai Yang, Xingzhong Xu, Xin Geng, and Xu Yang. Lmm-r1: Empowering 3b lmms with strong reasoning abilities through two-stage rule-based rl. *arXiv preprint arXiv:2503.07536*, 2025.

Yujia Qin, Yining Ye, Junjie Fang, Haoming Wang, Shihao Liang, Shizuo Tian, Junda Zhang, Jiahao Li, Yunxin Li, Shijue Huang, et al. Ui-tars: Pioneering automated gui interaction with native agents. *arXiv preprint arXiv:2501.12326*, 2025.

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, 2023.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.

Haozhan Shen, Zilun Zhang, Kangjia Zhao, Qianqian Zhang, Ruochen Xu, and Tiancheng Zhao. Vlm-r1: A stable and generalizable r1-style large vision-language model. https://github.com/om-ai-lab/VLM-R1, 2025. Accessed: 2025-02-15.

Junyang Wang, Haiyang Xu, Haitao Jia, Xi Zhang, Ming Yan, Weizhou Shen, Ji Zhang, Fei Huang, and Jitao Sang. Mobile-agent-v2: Mobile device operation assistant with effective navigation via multi-agent collaboration. *arXiv preprint arXiv:2406.01014*, 2024a.

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.

Weiyun Wang, Zhangwei Gao, Lianjie Chen, Zhe Chen, Jinguo Zhu, Xiangyu Zhao, Yangzhou Liu, Yue Cao, Shenglong Ye, Xizhou Zhu, et al. Visualprm: An effective process reward model for multimodal reasoning. *arXiv preprint arXiv:2503.10291*, 2025.

Zhiyong Wu, Zhenyu Wu, Fangzhi Xu, Yian Wang, Qiushi Sun, Chengyou Jia, Kanzhi Cheng, Zichen Ding, Liheng Chen, Paul Pu Liang, et al. Os-atlas: A foundation action model for generalist gui agents. *arXiv preprint arXiv:2410.23218*, 2024.

Zongru Wu, Pengzhou Cheng, Zheng Wu, Tianjie Ju, Zhuosheng Zhang, and Gongshen Liu. Smoothing grounding and reasoning for mllm-powered gui agents with query-oriented pivot tasks, 2025. URL https://arxiv.org/abs/2503.00401.

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. *arXiv preprint arXiv:2412.04454*, 2024.

Weihao Zeng, Yuzhen Huang, Wei Liu, Keqing He, Qian Liu, Zejun Ma, and Junxian He. 7b model and 8k examples: Emerging reasoning with reinforcement learning is both effective and efficient. https://hkust-nlp.notion.site/simplerl-reason, 2025. Notion Blog.

Chi Zhang, Zhao Yang, Jiaxuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. Appagent: Multimodal agents as smartphone users. *arXiv preprint arXiv:2312.13771*, 2023.

# A  Training

## A.1  Setting

We configure the hyperparameters as listed in Table 4 and train the base model using 8 NVIDIA 4090 GPUs, completing the training process in approximately 8 hours.

| Hyperparameter | Value |
|---|---|
| lr | from 9.98e-7 to 0 |
| max_pixels | 12845056 |
| num_generations | 8 |
| num_train_epochs | 8 |
| max_prompt_length | 1024 |
| per_device_train_batch_size | 1 |
| gradient_accumulation_steps | 2 |

Table 4: Hyperparameter settings used in the experiments.

## A.2  Dataset Distribution

The distribution of our data selection is listed in Table 5.

| Trainng dataset | Type | # Click | # Scroll | # Input text | # Back | # Open app | # Total |
|---|---|---|---|---|---|---|---|
| UI-R1 | Mobile | 101 | 6 | 4 | 7 | 18 | 136 |
| **Evaluation dataset** | | | | | | | |
| AndroidControl | ID | 5074 | 1211 | 632 | 343 | 608 | 7868 |
| ScreenSpot* | OOD | 770 | 0 | 0 | 0 | 0 | 770 |
| ScreenSpot-pro | OOD | 1581 | 0 | 0 | 0 | 0 | 1581 |

Table 5: Statistics of training and evaluation datasets. * means that we only select subsets Desktop and Web for evaluation.

## A.3 Visualization

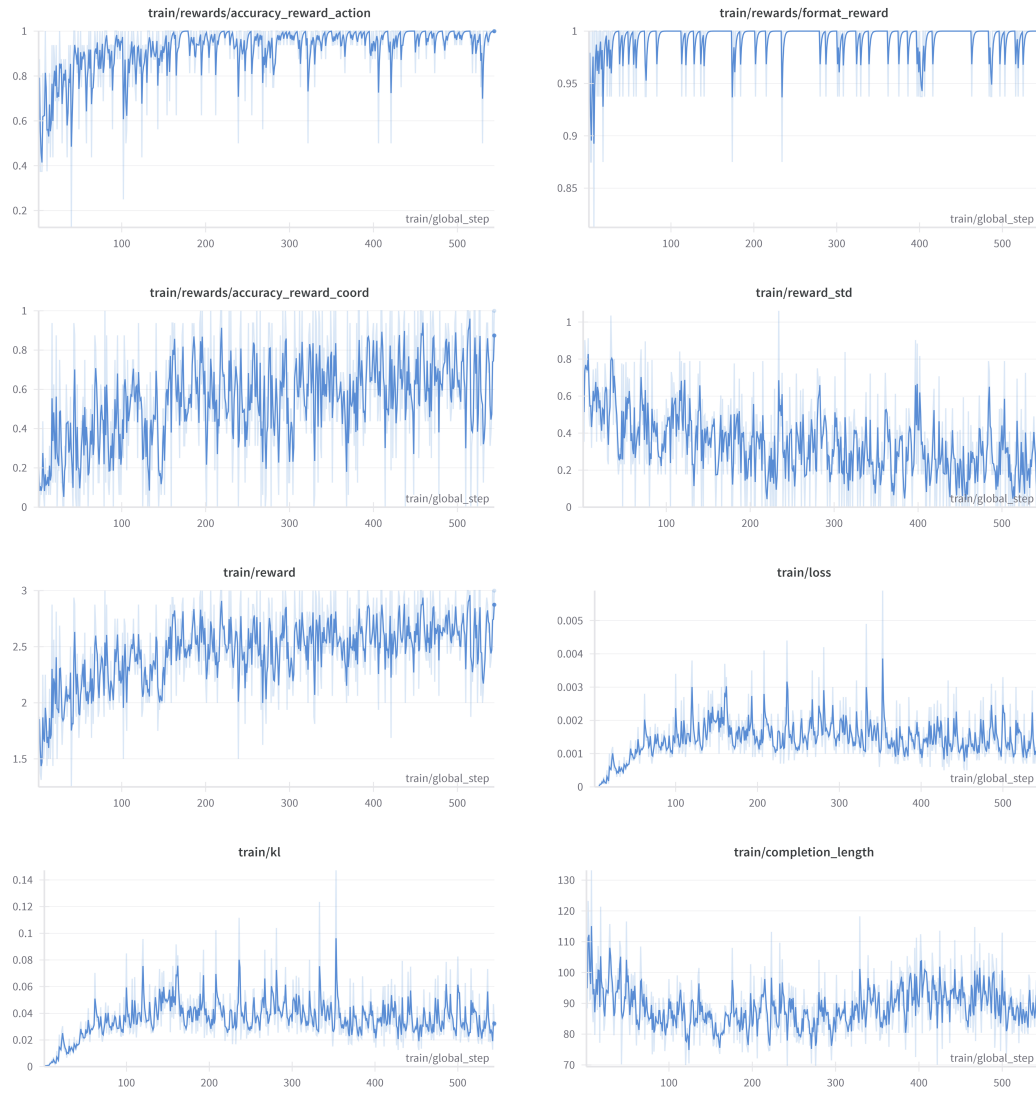Figure 5 illustrates the progression of various variables throughout the training process.



Figure 5: UI-R1 training process.

# B  Other Evaluation

## B.1  ScreenSpot-V2

We also evaluate the model performance on ScreenSpot-V2 (Wu et al., 2024) and the results are in Table 6.

| Model | GUI specific | Size | Mobile | | Web | | Desktop | | Avg |
|---|---|---|---|---|---|---|---|---|---|
| | | | Icon | Text | Icon | Text | Icon | Text | |
| SeeClick | Yes | 9.6B | 50.7 | 78.4 | 32.5 | 55.2 | 29.3 | 70.1 | 55.5 |
| OS-Atlas | Yes | 4B | 59.7 | 87.2 | 63.1 | 85.9 | 46.4 | 72.7 | 71.9 |
| OS-Atlas | Yes | 7B | 75.8 | 95.2 | <u>77.3</u> | **90.6** | <u>63.6</u> | <u>90.7</u> | 84.1 |
| UI-TARS | Yes | **2B** | 79.1 | 95.2 | **78.3** | 87.2 | **68.6** | <u>90.7</u> | <u>84.7</u> |
| **Qwen2.5-VL Framework** | | | | | | | | | |
| Qwen2.5-VL | No | 3B | 66.8 | 92.1 | 46.8 | 72.6 | 44.3 | 83.0 | 70.4 |
| Qwen2.5-VL | No | 7B | <u>80.6</u> | <u>95.9</u> | 70.0 | 87.2 | 59.3 | 89.2 | 82.6 |
| UI-R1(Ours) | Yes | <u>3B</u> | **84.3** | **96.2** | 75.4 | <u>89.2</u> | <u>63.6</u> | **92.3** | **85.4** |

Table 6: Grounding accuracy on ScreenSpot-V2. The optimal and the suboptimal results are **bolded** and <u>underlined</u>, respectively.

# C  Other Ablation

## C.1  Training epoches

We evaluate the model's performance across different training epochs, as shown in Figure 6. Based on the results, we finalize the training at 8 epochs.
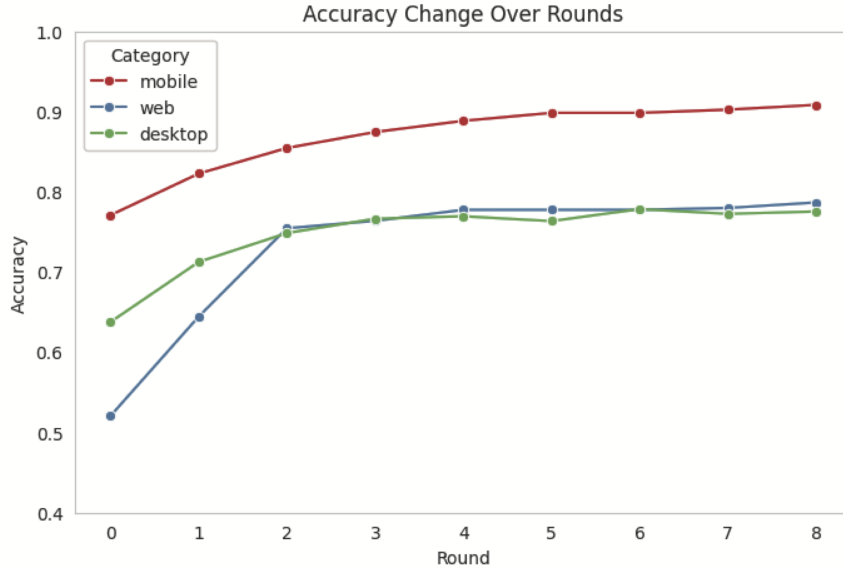


Figure 6: Accuracy change over rounds.

## C.2 Max Pixels

While adjusting the parameters, we observe that the maximum pixel setting of the image processor plays a significant role. If the input image exceeds this maximum pixel value, the `smart resize` function automatically crops and resizes the image while preserving the original aspect ratio. Mobile images are typically smaller than web or desktop images and often have significantly different aspect ratios. To address this, we implement the algorithm to appropriately rescale the predicted coordinates as shown in Algorithm 1.

---

**Algorithm 1** Scale Coordinates Based on Image Resizing

---

1: **Input:**
2:    $C = (x, y)$ : coordinate
3:    $I$ : input image
4:    $max\_pixels$ : maximum pixel constraint
5: **Output:** $(x_{scale}, y_{scale}) \in \mathbb{R}^2$
6: $(origin\_width, origin\_height) \leftarrow I.size$
7: $(resized\_height, resized\_width) \leftarrow \texttt{smart\_resize}(origin\_height, origin\_width, max\_pixels)$
                                                 ▷ `smart_resize` from QwenVL
8: $x_{scale} \leftarrow origin\_width / resized\_width \texttt{ * } x$
9: $y_{scale} \leftarrow origin\_height / resized\_height \texttt{ * } y$
10: **Output:** $(x_{scale}, y_{scale})$

---

We also investigate the impact of the maximum pixel value on model performance. Setting this value too high can lead to out-of-memory (OOM) errors during training when processing large images. Conversely, setting it too low may negatively affect the accuracy of prediction results. To better understand this trade-off, we experiment with two different maximum pixel values during training and evaluation, as summarized in Table 7.
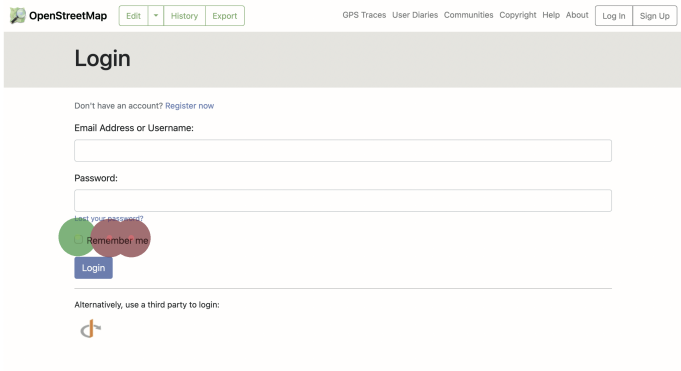
Based on our analysis, we set the maximum pixel value to 12,845,056 during training, which results in a model with improved performance on out-of-domain tasks. For evaluation, we recommend using a smaller maximum pixel value to conserve memory.

| max_pixels | | Mobile | Web | Desktop | Avg |
| Train | Test | | | | |
|---|---|---|---|---|---|
| 3211264 | 3211264 | **91.2** | 76.1 | 76.6 | 82.2 |
| 3211264 | 12845056 | 90.8 | 76.8 | 76.6 | 82.3 |
| 12845056 | 3211264 | 89.6 | 78.0 | **77.8** | 82.5 |
| 12845056 | 12845056 | 90.8 | **79.6** | 77.2 | **83.4** |

Table 7: Ablation of max pixels in the training and inference.

## D Case Study

Figure 7 illustrates an example of how UI-R1 trained model can successfully complete the task.



Figure 7: An example of use case.