

UI-AGILE: Advancing GUI Agents with Effective Reinforcement Learning and Precise Inference-Time Grounding

Shuquan Lian¹, Yuhang Wu¹, Jia Ma¹, Zihan Song¹, Bingqi Chen¹, Xiawu Zheng¹, Hui Li¹

¹Key Laboratory of Multimedia Trusted Perception and Efficient Computing Ministry of Education of China, Xiamen University

Abstract

The emergence of Multimodal Large Language Models (MLLMs) has driven significant advances in Graphical User Interface (GUI) agent capabilities. Nevertheless, existing GUI agent training and inference techniques still suffer from a dilemma for reasoning designs, ineffective reward, and visual noise. To address these issues, we introduce UI-AGILE, a comprehensive framework enhancing GUI agents at both the training and inference stages. For training, we propose a suite of improvements to the Supervised Fine-Tuning (SFT) process: 1) a Continuous Reward function to incentivize high-precision grounding; 2) a “Simple Thinking” reward to balance planning with speed and grounding accuracy; and 3) a Cropping-based Resampling strategy to mitigate the sparse reward problem and improve learning on complex tasks. For inference, we present Decomposed Grounding with Selection, a novel method that dramatically improves grounding accuracy on high-resolution displays by breaking the image into smaller, manageable parts. Experiments show that UI-AGILE achieves the state-of-the-art performance on two benchmarks ScreenSpot-Pro and ScreenSpot-v2. For instance, using both our proposed training and inference enhancement methods brings 23% grounding accuracy improvement over the best baseline on ScreenSpot-Pro.

1 Introduction

Driven by the growing capabilities of Multimodal Large Language Models (MLLMs) in image understanding (Bai et al. 2025), Graphical User Interface (GUI) agents, which execute tasks by understanding screenshots and user instructions, are advancing rapidly (Zhang et al. 2025).

Prior methods for GUI agents mostly rely on Supervised Fine-Tuning (SFT), requiring a large amount of human-annotated or synthesized data for teaching the agent how to plan its actions and grounding. (Qin et al. 2025; Cheng et al. 2024; Gou et al. 2025; Wu et al. 2025; Xu et al. 2024; Lin et al. 2025). Recently, Reinforcement Fine-Tuning (RFT) has emerged as an efficient and scalable way to instill reasoning abilities in MLLMs (Chen et al. 2025; Chen, Luo, and Li 2025; Liu et al. 2025b; Peng et al. 2025; Meng et al. 2025). For instance, UI-R1 (Lu et al. 2025) and GUI-R1 (Luo et al. 2025) apply RFT to enhance GUI agents and achieve considerable improvements compared to previous approaches.

Despite the significant momentum of GUI agent techniques, their practical application is hindered by several limitations in both training and inference stages:

- **P1: A Dilemma for Reasoning Designs:** An elaborate reasoning process not only degrades grounding accuracy but also increases inference latency, while a “No Thinking” approach exhibits low accuracy for predicting non-grounding action types (Shen et al. 2025).
- **P2: Ineffective Reward:** Agents often get stuck on complex interfaces and receive no effective learning signal (i.e., sparse reward). Besides, simple binary feedback (correct/incorrect), a design used by many existing methods (Lu et al. 2025; Luo et al. 2025) may fail to endow agents with the ability to perform precise localization.
- **P3: Visual Noise:** Even well-trained agents frequently struggle to cope with high-resolution screens, as irrelevant visual noise degrades their grounding accuracy.

To address the above problems, we propose UI-AGILE, a framework aiming at improving both the RFT and the inference stages of GUI agents. The contributions of UI-AGILE can be summarized as follows:

- To overcome **P1**, UI-AGILE applies a “Simple Thinking” strategy (Sec. 3.1) that employs reasoning through thoughts with appropriate lengths. UI-AGILE operationalizes “Simple Thinking” through a specialized reward function. “Simple Thinking” effectively balances the improvement of both the core grounding task and the prediction of non-grounding action types.
- To tackle **P2**, UI-AGILE harnesses a design of continuous grounding reward (Sec. 3.2) for the RFT stage instead of using the common binary reward to incentivize more precise localization to the target’s center. Furthermore, UI-AGILE employs cropping-based resampling (Sec. 3.3) to dynamically adjust the difficulty of training samples to avoid ineffective training with zero reward.
- To solve **P3**, UI-AGILE uses a visual noise reduction method termed decomposed grounding with selection (Sec. 3.4). It decomposes a high-resolution screenshot into multiple sub-images, generates candidate elements on each, and finally uses a Vision-Language Model (VLM) to “adjudicate” the best match. This approach significantly improves the agent’s grounding accuracy on high-resolution displays during inference.

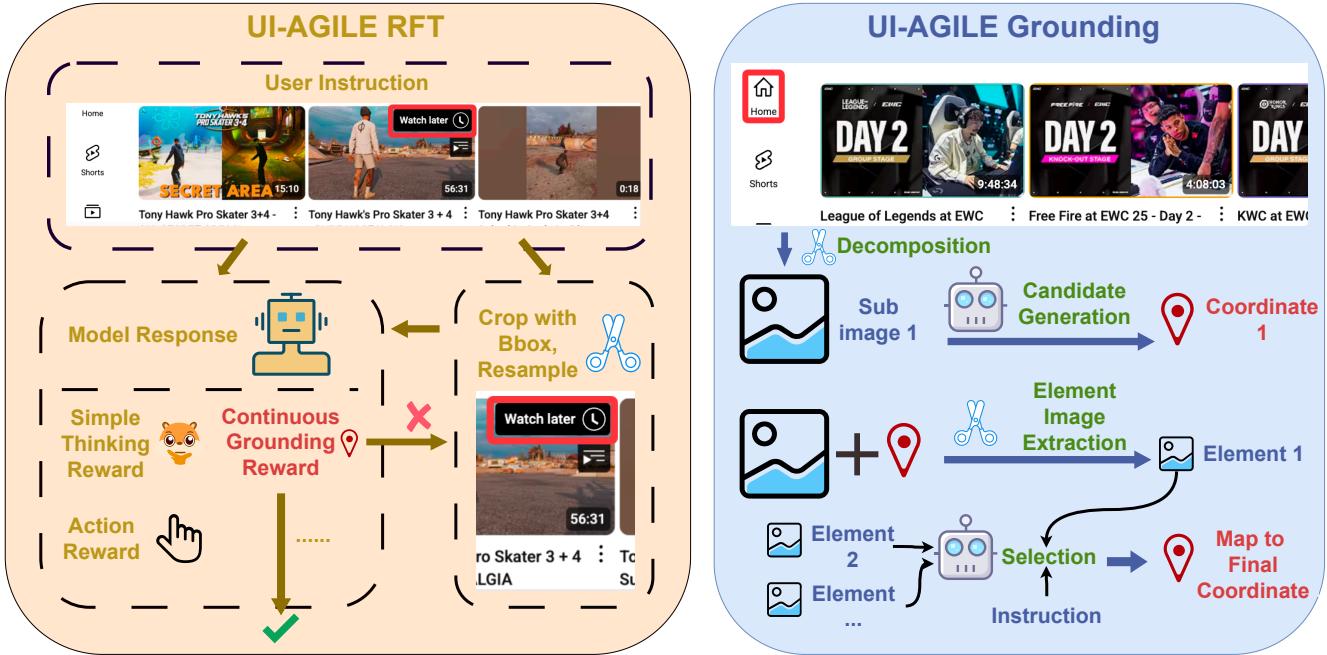


Figure 1: An overview of UI-AGILE, illustrating its training and inference pipelines. (1) Left: The training pipeline enhances the RFT process with our three core contributions: “Simple Thinking”, continuous grounding reward and cropping-based resampling. Continuous grounding reward being zero would result in crop-based resampling. (2) Right: The inference pipeline contains our proposed decomposed grounding with selection.

Extensive experiments validate the effectiveness of our methods. Trained on about only 9k samples for just 2 epochs, UI-AGILE shows superior performance, while also showcasing strong general agent capabilities. Furthermore, our inference method can act as a plug-and-play enhancement for a wide range of existing agents, improving the accuracy of some existing open-source models.

Overall, on two benchmarks ScreenSpot-Pro and ScreenSpot-v2, our methods achieve the state-of-the-art performance. For instance, using both our proposed training and inference enhancement methods brings 23% grounding accuracy improvement over the best baseline on ScreenSpot-Pro.

2 Related Work

Reinforcement Learning (RL) for Large Models Recently, reinforcement learning (RL) techniques for training large models have gained significant momentum. The focus has shifted from traditional policy optimization algorithms like PPO (Schulman et al. 2017), towards alignment-centric methods such as DPO (Rafailov et al. 2023), and more recent rule-based algorithms like GRPO (Shao et al. 2024). These algorithms have achieved remarkable success in enhancing the reasoning capabilities of large models on complex tasks, with models like OpenAI O1 (Jaech et al. 2024) and DeepSeek-R1 (DeepSeek-AI et al. 2025) setting new standards in mathematics and code generation. The efficacy of these approaches prompted their rapid extension into the multimodal domain (Chen et al. 2025; Chen, Luo, and Li

2025; Liu et al. 2025b; Peng et al. 2025; Meng et al. 2025).

GUI Agents The field of GUI agents has evolved rapidly (Wang et al. 2024; Li and Huang 2025). Following early works like CogAgent (Hong et al. 2024) and SeeClick (Cheng et al. 2024), recent studies rely on SFT to operate directly on visual inputs. Show-UI (Lin et al. 2025) innovates on visual processing efficiency. OS-Atlas (Wu et al. 2025), UGround (Gou et al. 2025) and Aria-UI (Yang et al. 2024) propose novel, large-scale pipelines to collect and synthesize millions of GUI agent trajectories, significantly improving model generalization. Aguvis (Xu et al. 2024) introduced a two-stage training process that explicitly uses VLM-generated Chain-of-Thought (CoT) data to teach planning and reasoning. JEDI (Xie et al. 2025) constructs a refusal part by mismatching existing instructions with unrelated screenshots. Standing out in complexity and scale, UI-TARS (Qin et al. 2025) utilizes the largest dataset and the most intricate training pipeline, which involves SFT and DPO on human annotated CoT data to improve performance. This data-intensive scaling has motivated a shift towards more efficient Reinforcement Learning (RL) paradigms, first explored by UI-R1 (Lu et al. 2025) and GUI-R1 (Luo et al. 2025). InfiGUI-R1 (Liu et al. 2025a) employs Spatial Reasoning Distillation and RL to enhance planning and error recovery capabilities. GUI-G1 (Zhou et al. 2025) leverages Hit-based reward and IoU-based reward for improving GUI agents.

3 Our Framework UI-AGILE

Fig. 1 provides an overview of UI-AGILE, which aims at improving both training and inference stages of GUI agents.

At training time, UI-AGILE adopts a new design for the reward function, which consists of a “Simple Thinking” reward for efficient reasoning (Sec. 3.1) and a continuous grounding reward (Sec. 3.2) for precise localization. UI-AGILE further uses cropping-based resampling (Sec. 3.3), a novel strategy designed to overcome the sparse reward issue. The model generates multiple responses from an image and instruction, which are evaluated by reward functions, including “Simple Thinking” reward and continuous grounding reward. If a training sample proves too difficult (i.e., receive a grounding reward of zero for all generated responses), the image will be cropped to simplify the task, and the model resamples new responses from this modified input.

At inference time, UI-AGILE applies decomposed grounding with selection (Sec. 3.4) to enhance grounding on the high-resolution displays common in modern applications, making GUI agents more practical for real-world use.

3.1 “Simple Thinking” for Reconciling the Reasoning Dilemma (P1)

GRPO algorithm (Shao et al. 2024) is widely recognized as a powerful technique for instilling complex reasoning abilities in LLMs, often through rewarding elaborate chains of thought. Lu et al. (2025) posit that excessive reasoning is not essential for GUI grounding and can even be detrimental. However, the complete role of GUI agents extends beyond mere grounding, since deciding next action (e.g., click or type) inherently demands a foundational level of reasoning.

To reconcile the dilemma of whether to apply reasoning for enhancing action-type prediction or discard excessive reasoning to avoid ineffective grounding, we propose a “Simple Thinking” strategy. It encourages thoughts with an appropriate length, operationalized through a specialized reward function. It also reduces training and inference costs, a practical consideration for real-world deployment.

The total reward for the thought process, R_{think} is formally defined as:

$$R_{\text{think}} = I(R_{\text{grounding}} > 0) \cdot (R_{\text{length}}(L) + R_{\text{bonus}}) \quad (1)$$

$$R_{\text{length}}(L) = \begin{cases} 1.0 & \text{if } l_{\text{ideal.start}} < L \leq l_{\text{ideal.end}} \\ \frac{1}{2} \left(1 - \cos \left(\pi \frac{L - l_{\text{min}}}{l_{\text{ideal.start}} - l_{\text{min}}} \right) \right) & \text{if } l_{\text{min}} < L \leq l_{\text{ideal.start}} \\ \frac{1}{2} \left(1 + \cos \left(\pi \frac{L - l_{\text{ideal.end}}}{l_{\text{max}} - l_{\text{ideal.end}}} \right) \right) & \text{if } l_{\text{ideal.end}} < L < l_{\text{max}} \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

where:

- $I(R_{\text{grounding}} > 0)$ is an indicator function that grants reward only when the grounding reward $R_{\text{grounding}} > 0$, linking reasoning to effective outcomes.
- $R_{\text{length}}(L)$ is a non-linear reward based on the reasoning length L . $l_{\text{ideal.start}}$ and $l_{\text{ideal.end}}$ define an ideal range of reasoning length, where reward is 1. The reward will be zero if the reasoning length exceeds l_{min} or l_{max} .

- R_{bonus} is a fixed bonus for syntactically complete thoughts (e.g., ending with proper punctuation), encouraging structured reasoning.

“Simple Thinking” defines an ideal range where the reward is maximized, encouraging thoughts that are neither too brief (“under-thinking”) nor too verbose (“over-thinking”). Outside this ideal range, it uses cosine functions for smooth degradation down to a reward of zero at the absolute bounds. This smooth, non-linear penalty provides a more informative and stable learning signal for RL than a hard cliff, gently guiding the model toward the desired length. Furthermore, the additional bonus for syntactically complete thoughts discourages incomplete reasoning, thereby ensuring better training stability.

3.2 Continuous Grounding Reward for Precise Localization (P2)

Prior works (Lu et al. 2025; Luo et al. 2025) shift the focus of GUI agent evaluation from traditional object grounding (i.e., IoU) to the precision of the action coordinate. To this end, they typically employ a simple binary reward for localization: a reward of 1 for a correct prediction (e.g., inside the target radius) and 0 otherwise.

However, this binary reward is insufficient for high-precision control, as it treats all successful predictions equally. For instance, a point on the target’s edge receives the same reward as one at its center. This non-discriminatory feedback misguides the model to learn an element’s boundaries rather than its semantic core.

To resolve this issue, we introduce a continuous grounding reward. Instead of a binary outcome, the continuous reward is calculated as a function of the distance from the predicted point to the center of the ground-truth bounding box:

$$R(x, y) = \begin{cases} 1 + \exp(-4 \cdot d_{\text{norm}}^2) & \text{if } (x, y) \in \text{BBox} \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

where:

- $R(x, y)$ is the reward score for the predicted coordinate.
- (x, y) is the coordinate predicted by the agent.
- BBox is the ground-truth bounding box, defined by its top-left (x_1, y_1) and bottom-right (x_2, y_2) corners.
- d_{norm} is the Chebyshev distance (or L_∞ norm) of the point from the center of the bounding box, normalized by the box’s dimensions. It is calculated as:

$$d_{\text{norm}} = \max \left(\frac{|x - c_x|}{w_h}, \frac{|y - c_y|}{h_h} \right) \quad (4)$$

Here, $(c_x, c_y) = (\frac{x_1+x_2}{2}, \frac{y_1+y_2}{2})$ represents the center of the bounding box, and $(w_h, h_h) = (\frac{x_2-x_1}{2}, \frac{y_2-y_1}{2})$ are its half-width and half-height, respectively.

We employ the Chebyshev distance instead of the Euclidean distance because the reward contours generated by the Chebyshev distance are squares, which geometrically align with the rectangular shape of GUI bounding boxes. This alignment provides a more logical and consistent reward landscape, strongly incentivizing the agent to minimize its deviation along both axes to achieve a high reward.

3.3 Cropping-Based Resampling for Sparse Reward Mitigation (P2)

During GRPO training, GUI agents often face the sparse reward challenge, particularly on complex tasks. When the model consistently fails to place its prediction within the correct bounding box for a given screenshot, it receives a persistent reward of 0. The lack of positive signals can lead to training stagnation, and difficult samples cannot contribute to model improvement.

Inspired by curriculum learning (Bengio et al. 2009), a training technique that trains the model on examples of increasing difficulty to ease the training process, we propose a cropping-based resampling method for mitigating the sparse reward issue. It acts as a dynamic difficulty adjustment mechanism during the training of UI-AGILE. If a task sample yields zero reward over multiple generations, we hypothesize that the task sample is currently too difficult for the model. Then, we reduce the task’s complexity by cropping the original screenshot. The cropping is generated such that it is smaller than the original view but still fully contains the ground-truth bounding box of the target element.

A naive implementation is to center the ground-truth bounding box (bbox) in the new cropping, but it may fail to endow the model with the capability to perform robust localization: the model would learn a trivial shortcut, such as developing a bias for predicting the image center. We opt to employ a scanning approach as illustrated in Alg. 1, ensuring that the cropped image fully contains the ground-truth bounding box. It firstly determines the size of cropping based on a predefined ratio (lines 1-2). Then, the horizontal stride $step_x$ is set to the difference between the cropping width and the bounding box width, while the vertical stride $step_y$ is set to the difference between their respective heights (lines 3-5). After that, it iterates through all possible cropping windows from left-to-right and top-to-bottom across the original screenshot with the horizontal stride and the vertical stride (lines 7-18). The first window that fully contains the ground-truth bbox is selected as the new, resampled input (lines 12-16). Fig. 2 illustrates how our scanning approach identifies valid cropping windows that fully contain the ground-truth bounding box.

Cropping-based resampling dynamically simplifies difficult samples to ensure that they are learnable, allowing the model to leverage more data in fewer epochs, yielding superior results within a similar amount of training time.

3.4 Decomposed Grounding with Selection for Visual Noise Reduction (P3)

Modern electronic devices feature high-resolution displays (e.g., 3840x2160), which, when converted into tokens for a VLM, can result in an overwhelmingly large input sequence (e.g., over 10,000 tokens). We hypothesize that a significant portion of these tokens represent irrelevant background information, acting as noise that can distract the GUI agent and degrade its grounding accuracy, i.e., the visual noise issue (**P1**) mentioned in Sec. 1.

To validate this hypothesis, we conduct a preliminary experiment on ScreenSpot-Pro (Li et al. 2025). We apply the

Algorithm 1: Cropping-Based Resampling

```

Input: Image, Bbox, scalingFactorf
Output: Set(Imagecrop, Bboxcrop)
1: (wo, ho)  $\leftarrow$  GetWidthAndHeight(Image)
2: (wcrop, hcrop)  $\leftarrow$  (wo  $\times$  f, ho  $\times$  f)
3: (wb, hb)  $\leftarrow$  (Bbox[2]  $-$  Bbox[0], Bbox[3]  $-$  Bbox[1])
4: stepx  $\leftarrow$  wcrop  $-$  wb
5: stepy  $\leftarrow$  hcrop  $-$  hb
6: T  $\leftarrow$   $\emptyset$ 
7: for xcropmin from 0 to wo step=stepx do
8:   xcropmax  $\leftarrow$  min(xcropmin + wcrop, wo)
9:   for ycropmin from 0 to ho step=stepy do
10:    ycropmin  $\leftarrow$  min(ycropmin + hcrop, ho)
11:    Coordcrop  $\leftarrow$  [xcropmin,  
           ycropmin, xcropmax, ycropmax]
12:    if Bbox is contained within Coordcrop then
13:      Imagecrop  $\leftarrow$  Crop(Image, Coordcrop)
14:      Bboxcrop  $\leftarrow$  Bbox  $-$  Coordcrop
15:      Set  $\leftarrow$  Set  $\cup$  (Imagecrop, Bboxcrop)
16:    end if
17:  end for
18: end for
19: return Set

```

cropping method in Sec. 3.3, but for the purpose of creating a controlled test environment. For each original screenshot, we crop it to 1024x1024, ensuring the ground-truth bounding box is contained within the frame. On this new dataset, the grounding accuracy of UGround-V1-7B (Gou et al. 2025) shows a significant improvement from 31.6 to 56.0 w.r.t grounding accuracy, verifying our hypothesis.

It is straightforward to consider applying the cropping-based method to alleviate the visual noise. However, during inference, the location of the ground-truth bounding box is unknown, making such oracle-based cropping impossible. To overcome this challenge, we propose a multi-stage decomposed grounding with selection method as shown in the right part of Fig. 1. It reduces visual noise by cropping an image into several sub-images and predicts the coordinate individually, while trying to maintain full ground-truth bounding box area. The detailed process is as follows:

1. **Decomposition:** The input screenshot is first divided into several overlapping sub-images, breaking down the high-resolution screen into smaller, more manageable regions.
2. **Candidate Generation:** The GUI agent performs grounding independently on each sub-image and predicts a coordinate for the sub-image. Then, we use the predicted coordinates of sub-images belonging to a screenshot as the candidate coordinates for the screenshot.
3. **Element Image Extraction:** For each candidate point, we extract the corresponding element image by cropping a bounding box centered on the candidate point from the sub-image.
4. **Selection:** For the final selection stage, we prompt the VLM with the user’s instruction, the candidate element image, and a direct question asking whether the image match the instruction. Then use the model’s output logit for the “Yes” token as a direct relevance score for the candidate. The candidate with the highest “Yes” score

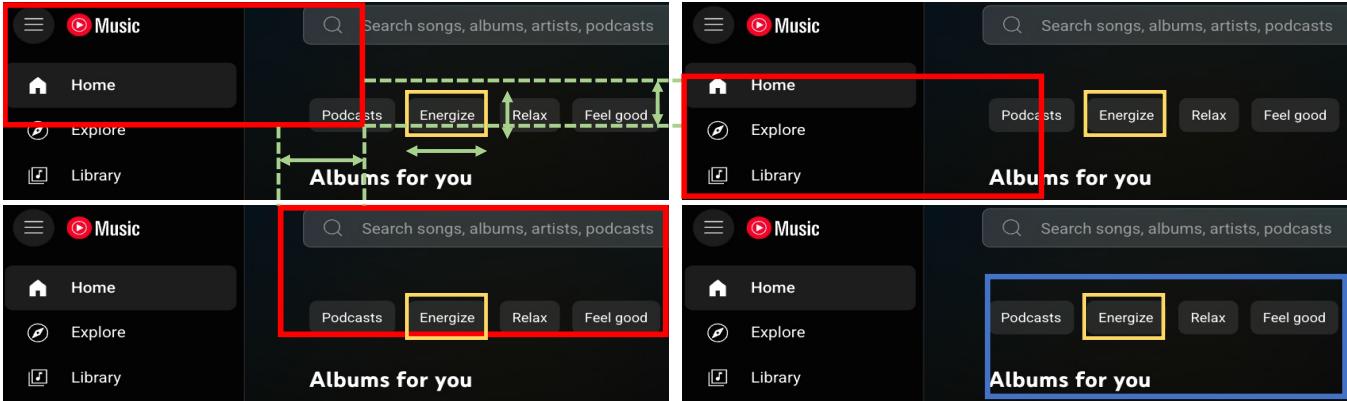


Figure 2: An example of cropping-based resampling. Yellow bounding boxes are the ground truth, red bounding boxes are invalid cropping, and blue bounding boxes are valid cropping. Green arrows show that the overlap of cropping windows is equivalent to the width or height of the ground-truth bounding box.

is chosen as the final answer, and its corresponding coordinates are remapped to the original screenshot. This QA-based scoring allows the model to perform deeper contextual prediction.

This process directly benefits from our continuous reward function (Sec. 3.2), as it trains the model to predict points closer to the center of target elements, leading to higher-quality extracted images and thus a more accurate final selection.

Analysis of Inference Cost We analyze the inference latency of using decomposed grounding with selection by breaking it down into three primary stages: prefilling, decoding and selection stage.

Counter-intuitively, our approach can *theoretically accelerate the computationally-heavy prefilling stage*. Recall that the self-attention mechanism has a quadratic time complexity $\mathcal{O}(n^2)$ concerning the input sequence length n (Vaswani et al. 2017). By splitting a single large image with n tokens into 4 sub-images (each with roughly $n/4$ tokens), the total computational cost for attention scales proportionally to $4 \times (\frac{n}{4})^2 = \frac{n^2}{4}$, suggesting a theoretical speedup. Such cost reduction far outweighs the slightly increased overhead of repeatedly processing the text prompt for each sub-image.

While the decoding process runs for each sub-image, the small number of output tokens (compared to image tokens) per run ensures that the cumulative decoding cost will not raise much.

Finally, the VLM-based selection stage is computationally inexpensive. The input element images are very small, and the process only requires a single forward pass to acquire the logits for a “Yes/No” answer.

Overall, the cost of applying decomposed grounding with selection is low. In addition to the above analysis, we provide results on actual running time in Sec. 4.6. Crucially, we believe this overhead could be eliminated or even reversed with future optimizations in inference engines tailored for this “many small requests” workload, potentially making our high-accuracy method faster than the baseline.

4 Experiment

4.1 Implementation Details

Data. We collect data related to GUI tasks from multiple open-source datasets, including UI-R1 (Lu et al. 2025), GUI-R1 (Luo et al. 2025), Aguvix (Xu et al. 2024) and Grounding-R1 (Yang et al. 2025). We filter them using OmniParser (Wan et al. 2024) following Grounding-R1 (Yang et al. 2025). We randomly sample about 9k examples to train UI-AGILE-3B and UI-AGILE-7B.

Baselines. We include detailed descriptions for all baselines in Appendix.

Training Details. We use the trl framework (von Werra et al. 2020) to implement the cropping-based resampling strategy and reward functions. The sampling process is attempted 4 times at most and is bypassed entirely if the bbox’s dimensions exceed the target crop size. Following previous works (Lu et al. 2025; Luo et al. 2025; Liu et al. 2025a; Zhou et al. 2025), we use Qwen2.5-VL¹ as base model. Tab. 1 provides the default hyperparameters where cropping factor is the width and height ratio of new attempted image and last attempted image.

Hyperparameter	Value
learning rate	from 1e-06 to 4.36e-10
num generations	8
num train epochs	2
per device train batch size	4
cropping factor	0.6
gradient accumulation steps	4

Table 1: Hyperparameters

¹<https://huggingface.co/Qwen/Qwen2.5-VL-3B-Instruct>, <https://huggingface.co/Qwen/Qwen2.5-VL-7B-Instruct>

Inference Details. For decomposed grounding with selection, the input image is divided into four sub-images scaling to 60% of the original dimensions, with adjacent sub-images overlapping by 10% of the original image’s width and height. In the element image extraction stage, we define the element’s area by creating a simple bounding box centered on the predicted point with the width and height equal to 14% of the sub-image’s width and height. We have also explored a more sophisticated approach using OmniParser to refine this bounding box. However, it does not improve performance and increases the inference overhead. In the selection stage, we use Qwen2.5VL-7B-instruct¹ to choose the final answer.

4.2 Grounding Capability Evaluation

We evaluate the grounding ability on ScreenSpot-v2 (Wu et al. 2025) and ScreenSpot-Pro (Li et al. 2025). ScreenSpot-v2 is a corrected version of the original ScreenSpot (Cheng et al. 2024), providing evaluation of GUI grounding capability across mobile, desktop, and web platforms. ScreenSpot-Pro focuses on high-resolution professional environments, featuring expert-annotated tasks spanning 23 applications, five industries, and three operating systems. Since the images in ScreenSpot-v2 are already pre-cropped while those in ScreenSpot-Pro are full, uncropped displays, we evaluate our decomposed grounding with selection method exclusively on ScreenSpot-Pro.

Effectiveness of the proposed inference enhancement. As shown in Tab. 2, our decomposed grounding with selection method shows significant improvements on the challenging ScreenSpot-Pro benchmark. It provides a universal and substantial performance boost across all tested models, regardless of their original training paradigm (SFT or RFT). For instance, it elevates the average score of OS-Atlas-7B from 18.9 to 33.1 (**+75.1%**), and boosts Aguvis-7B from 20.4 to 36.5 (**+78.9%**). The consistent improvement validates the effectiveness of decomposed grounding with selection and its high applicability as a plug-and-play inference enhancement.

Effectiveness of the proposed training enhancement. Besides, Tab. 2 shows that our core UI-AGILE-3B and UI-AGILE-7B models, even without decomposed grounding, establish a new state-of-the-art among 3B and 7B models on ScreenSpot-Pro. Trained on only 9K examples with 2 epochs, they (37.9 for 3B and 44.0 for 7B) surpass other RFT-based models like UI-R1-E (33.5), InfiGUI-R1-3B and GUI-R1-7B (32.1). UI-AGILE-7B even outperforms the much larger model UI-TARS-72B (38.1) trained on approximately 50 billion tokens. On the ScreenSpot-v2 benchmark (Tab. 3), our UI-AGILE-7B also achieves state-of-the-art grounding accuracy with an average score of 92.1. The above results demonstrate the effectiveness of our proposed “Simple Thinking” reward, continuous grounding reward, and cropping-based resampling for improving the training of GUI agents.

Overall, as shown in Tab. 2, using both our proposed training and inference enhancement methods (UI-AGILE-7B + Decomposed Grounding) brings 23% grounding ac-

curacy improvement over the best baseline (JEDI-7B) on ScreenSpot-Pro.

4.3 Agent Capability Evaluation

In addition to grounding-specific benchmarks, we also evaluate UI-AGILE-3B and UI-AGILE-7B on AndroidControl (Li et al. 2024) to assess its general agent capabilities.

Following the evaluation setting of OS-Atlas (Wu et al. 2025), we use three metrics: action type prediction accuracy (Type), grounding accuracy (GR), and the overall step success rate (SR). Type accuracy measures the exact match for the predicted action (e.g., click or scroll). For GR, a prediction is considered successful if it falls within a 14% screen-width radius of the ground-truth coordinate. The most holistic metric, SR, deems a step successful only if both the action type and all its associated arguments (e.g., coordinates for a click, direction for a scroll, or text for an input) are correct. For the number of test samples on AndroidControl, we follow OS-Atlas and use 7,708 examples for a fair comparison, while some works (e.g., Aguvis (Xu et al. 2024) and UGround (Gou et al. 2025)) randomly sample 500 action steps from the full AndroidControl test set for testing.

The evaluation is conducted under two distinct settings. In the AndroidControl-Low setting, the agent receives a specific, low-level instruction for each step. In contrast, the AndroidControl-High setting provides the agent with a high-level goal, requiring it to infer the correct action for the current step based on the conversation history and previous actions, thereby testing its multi-step reasoning capability.

As presented in Tab. 4, our UI-AGILE-7B model achieves the best performance (SR: 77.6 and 60.6) compared to all other RFT models including UI-R1-E (SR: 71.37 and 35.88), GUI-R1-3B (SR: 64.41 and 46.55) and GUI-R1-7B (66.52 and 51.67). Notably, our smaller UI-AGILE-3B model with SR of 77.4 and 56.9 also surpasses 7B baselines like GUI-R1-7B (SR: 66.5 and 51.7). This remarkable result demonstrates that the improvements gained from our methods are not confined to improving grounding capability but also translate effectively to better decision-making and execution in complex, multi-step agent scenarios.

4.4 Ablation Study

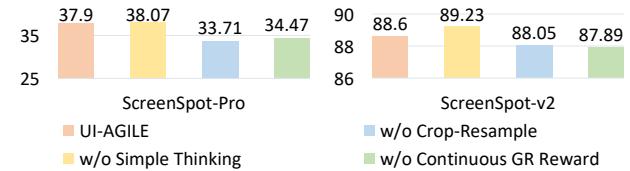


Figure 3: Ablation study on grounding benchmarks.

To further verify the contribution of each of our training techniques (“Simple Thinking” reward, continuous grounding reward, and cropping-based resampling), we conduct an ablation study of UI-AGILE-3B and the results are shown in Fig. 3 and Fig. 4. We can observe that:

Model	Examples	Epochs	Dev		Creative		CAD		Scientific		Office		OS		Avg	pass@4
			Text	Icon												
Supervised Fine-tuning																
CogAgent-18B	222M	-	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-UI	16.6M	-	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	-
ShowUI-2B	256K	-	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	-
JEDI-3B	4M	-	61.0	13.8	53.5	8.4	27.4	9.4	54.2	18.2	64.4	32.1	38.3	9.0	36.1	-
JEDI-7B	4M	-	42.9	11.0	50.0	<u>11.9</u>	38.0	14.1	72.9	25.5	<u>75.1</u>	47.2	33.6	16.9	39.5	-
OS-Atlas-7B	13M	-	33.8	1.4	30.8	3.5	12.2	3.1	33.3	9.1	33.3	9.4	26.2	3.4	18.9	-
+Decomposed Grounding	-	-	49.4	5.5	52.0	5.6	26.4	6.3	54.9	18.2	57.6	18.9	49.5	9.0	33.1	42.1
Aguvis-7B	1M	1	30.5	0.7	28.8	2.8	14.7	1.6	45.8	8.2	38.4	11.3	30.8	2.3	20.4	-
+Decomposed Grounding	-	-	50.6	11.7	60.1	7.0	31.0	4.7	62.5	20.0	63.3	18.9	45.8	6.7	36.5	44.3
UGround-V1-7B	10M	-	51.3	5.5	48.5	8.3	18.8	1.6	59.7	14.6	59.9	17.0	40.2	7.9	31.6	-
+Decomposed Grounding	-	-	57.8	14.5	49.0	<u>11.9</u>	20.3	7.8	62.5	21.8	67.8	18.9	48.6	14.6	36.6	47.3
UI-TARS-2B	-	-	47.7	4.1	42.9	6.3	17.8	4.7	56.9	17.3	50.3	17.0	21.5	5.6	27.7	-
UI-TARS-7B	-	-	58.4	12.4	50.0	9.1	20.8	9.4	63.9	31.8	63.3	20.8	30.8	16.9	35.7	-
+Decomposed Grounding	-	-	59.7	<u>19.3</u>	54.0	15.4	38.1	12.5	63.2	27.3	71.8	28.3	45.8	<u>21.3</u>	41.9	50.3
UI-TARS-72B	-	-	63.0	17.3	57.1	15.4	18.8	<u>17.2</u>	64.6	20.9	63.3	26.4	42.1	<u>15.7</u>	38.1	-
Zero Shot / Reinforcement Fine-tuning																
InfiGUI-R1-3B	32K	-	51.3	12.4	44.9	7.0	33.0	14.1	58.3	20.0	65.5	28.3	43.9	12.4	35.7	-
GUI-G1-3B	17K	1	50.7	10.3	36.6	<u>11.9</u>	39.6	9.4	61.8	<u>30.0</u>	67.2	32.1	32.5	10.6	37.1	-
Qwen2.5-VL-3B	-	-	31.8	4.1	32.8	4.2	24.9	4.7	43.8	12.7	42.4	15.1	17.8	2.2	22.7	-
+Decomposed Grounding	-	-	52.6	8.3	42.9	11.2	25.9	3.1	47.9	10.0	55.9	17.0	46.7	9.0	31.2	38.8
Qwen2.5-VL-7B	-	-	54.5	5.5	24.7	4.2	13.7	3.1	46.5	7.3	50.8	11.3	29.9	10.1	24.5	-
+Decomposed Grounding	-	-	60.4	13.1	33.3	8.4	27.9	6.2	50.0	13.6	63.3	17.0	51.4	16.0	33.3	42.0
GUI-R1-3B	3K	9	40.9	4.8	47.8	2.8	27.9	6.3	65.3	19.1	58.2	18.9	29.0	2.2	30.9	-
+Decomposed Grounding	-	-	63.6	13.1	55.6	4.9	31.5	6.3	61.8	16.4	62.7	20.8	44.9	10.1	37.1	46.0
GUI-R1-7B	3K	9	57.1	8.3	37.9	8.4	28.4	6.3	54.9	10.9	59.9	13.2	41.1	13.5	32.1	-
+Decomposed Grounding	-	-	<u>66.9</u>	15.2	50.0	10.5	32.5	4.7	59.0	13.6	68.4	24.5	60.7	18.0	39.3	48.6
UI-R1-3B	136	8	22.7	4.1	27.3	3.5	11.2	6.3	42.4	11.8	32.2	11.3	13.1	4.5	17.8	-
UI-R1-E	2K	8	46.1	6.9	41.9	4.2	37.1	12.5	56.9	21.8	65.0	26.4	32.7	10.1	33.5	-
+Decomposed Grounding	-	-	63.6	17.2	59.6	10.0	43.7	6.3	66.0	21.8	68.4	<u>43.4</u>	56.1	19.1	43.3	52.6
UI-AGILE-3B	9K	2	53.2	9.0	50.5	8.4	44.2	20.3	62.5	22.7	65.5	22.6	35.5	12.4	37.9	-
+Decomposed Grounding	-	-	<u>66.9</u>	16.6	58.1	<u>11.9</u>	47.2	10.9	66.7	24.5	72.3	34.0	58.9	22.5	45.0	54.4
UI-AGILE-7B	9K	2	64.3	15.2	53.0	9.8	<u>49.2</u>	14.1	72.9	25.5	<u>75.1</u>	30.2	45.8	20.2	44.0	-
+Decomposed Grounding	-	-	79.1	24.1	60.6	11.2	<u>53.3</u>	10.9	66.0	26.4	79.1	39.6	<u>59.8</u>	22.5	48.7	59.2

Table 2: Grounding accuracy on ScreenSpot-Pro. “+Decomposed Grounding” denotes that the model uses our decomposed grounding with selection for enhancing inference. Results marked in bold represent the best performance, and underlined results indicate the second-best performance. The pass@4 metric indicates the success rate where a task is considered solved if the prediction of at least one of the sub-images is correct.



Figure 4: Ablation study on AndroidControl benchmark.

- Applying continuous grounding reward and cropping-based resampling improve the performance by approximately 10% and 12.4% on ScreenSpot-Pro, respectively.

The former incentivizes more precise localization to the target’s center and the later helps avoid ineffective training with zero reward. They also slightly improve grounding accuracy on ScreenSpot-v2 where the performance of the base model is already high and it is difficult to achieve significant gains.

- Removing “Simple Thinking” during training (i.e., “No Thinking”) leads to higher grounding accuracy on ScreenSpot-Pro and ScreenSpot-v2 (about 0.4% and 0.7% gains). However, Fig. 4 demonstrates that integrating “Simple Thinking” enhances the agent’s decision-making with planning, resulting in a noticeable improvement.

Method	Mobile		Desktop		Web		Avg.
	Text	Icon	Text	Icon	Text	Icon	
SeeClick	78.4	50.7	70.1	29.3	55.2	32.5	55.1
OS-Atlas-7B	95.0	73.3	92.8	64.9	89.6	72.4	83.7
Aguvis-7B	94.9	80.1	95.0	77.9	91.4	69.9	85.6
Qwen2.5-VL-3B	96.1	74.8	87.8	53.0	86.9	70.4	80.7
Qwen2.5-VL-7B	98.4	84.8	88.4	74.7	92.5	77.6	87.5
GUI-R1-3B	98.1	79.0	94.0	66.7	93.3	69.2	85.2
GUI-R1-7B	98.8	86.4	92.3	79.4	92.1	77.2	88.7
UI-R1-E	<u>99.6</u>	80.1	95.6	75.4	91.6	81.2	88.7
UI-TARS-2B	95.2	79.1	90.7	68.6	87.2	78.3	84.7
UI-TARS-7B	96.9	<u>89.1</u>	<u>95.4</u>	85.0	93.6	85.2	<u>91.6</u>
UI-TARS-72B	94.8	86.3	91.2	87.9	91.5	87.7	90.3
UI-AGILE-3B	<u>99.6</u>	86.4	93.9	74.5	91.8	77.6	88.6
UI-AGILE-7B	100.0	91.1	95.6	84.8	94.2	83.0	92.1

Table 3: Grounding accuracy on ScreenSpot-v2. Results marked in bold represent the best performance, and underlined results indicate the second-best performance.

Models	AndroidControl-Low			AndroidControl-High		
	Type	GR	SR	Type	GR	SR
Os-Atlas-4B	64.6	71.2	40.6	49.0	49.5	22.8
Os-Atlas-7B	73.0	73.4	50.9	57.4	54.9	29.8
Qwen2.5-VL-3B	80.5	79.4	67.8	64.4	46.1	44.4
Qwen2.5-VL-7B	78.0	87.1	68.7	69.1	59.1	50.1
UI-R1-E	<u>87.0</u>	77.8	71.4	66.4	37.8	36.9
GUI-R1-3B	83.7	81.6	64.4	58.0	56.2	46.5
GUI-R1-7B	85.2	84.0	66.5	71.6	65.6	51.7
UI-AGILE-3B	85.4	<u>87.6</u>	<u>74.3</u>	<u>78.6</u>	60.7	56.9
UI-AGILE-7B	87.7	88.1	77.6	80.1	61.9	60.6

Table 4: Action type accuracy, grounding accuracy and success rate on AndroidControl-Low and AndroidControl-High. Results marked in bold represent the best performance, and underlined results indicate the second-best performance.

ment in SR on the AndroidControl benchmark, with SR increases of 15.5% and 3.4% in Low and High settings, respectively.

4.5 Analysis of Attempts Per Step

Figure 5 shows the distribution of attempts per GRPO training step, where each step processes a batch of two training samples. In the first epoch, we find that only 61.8% of training steps are fully successful on the initial attempt (i.e., both samples in the batch are solved without resampling). This means that without our strategy, a minimum of 19.1% ($38.2\% \div 2$) of training samples would have provided no learning signal. Overall attempt numbers decreases in the second epoch, demonstrating that the model learns from the samples salvaged by our method.

4.6 Analysis of Inference Time

We report the inference time of our decomposed grounding with selection method on the full ScreenSpot-Pro dataset (Li

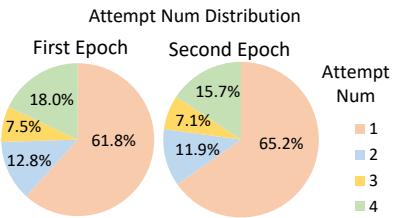


Figure 5: Distribution of attempts per step throughout the training process. Max attempt number is set to 4.

et al. 2025) using the vLLM framework (Kwon et al. 2023) and one 80G A800 GPU card. As a baseline, the standard grounding approach applied to UI-AGILE-7B without our method completes the benchmark in **30 minutes**. When applying our method, the decomposed grounding stage takes **35 minutes**. The subsequent VLM-based selection stage requires additional **4 minutes**. The modest increase in overhead is a practical trade-off for the substantial gain of grounding accuracy brought by our method.

5 Conclusions

In this paper, we introduce UI-AGILE, a comprehensive framework designed to enhance GUI agents’ training and inference. We tackle the practical challenges of the reasoning-grounding dilemma, ineffective reward, and visual noise. During training, our solution integrates three key innovations: a “Simple Thinking” reward to foster efficient yet effective reasoning, a continuous grounding reward that incentivizes high-precision localization, and a cropping-based resampling strategy to overcome the sparse reward problem. For inference, we introduce decomposed grounding with selection, a novel method that reduces visual noise and dramatically improves grounding accuracy on high-resolution screens while the inference cost is only slightly increased. Experimental results demonstrate the effectiveness of our proposed techniques on enhancing GUI agents’ grounding capability and general capability.

Limitations and future work. The VLM used in the selection stage of our decomposed grounding with selection method is a general-purpose, pre-trained model. A promising future direction would be to fine-tune this adjudicator model on a curated dataset of candidate UI elements, enhancing its selection accuracy to achieve further gains in overall grounding performance.

References

- Bai, S.; Chen, K.; Liu, X.; Wang, J.; and et al, W. G. 2025. Qwen2.5-VL Technical Report. *arXiv Preprint*.
- Bengio, Y.; Louradour, J.; Collobert, R.; and Weston, J. 2009. Curriculum learning. In *ICML*, volume 382 of *ACM International Conference Proceeding Series*, 41–48.
- Chen, L.; Li, L.; Zhao, H.; Song, Y.; and Vinci. 2025. R1-V: Reinforcing Super Generalization Ability in Vision-Language Models with Less Than \$3. <https://github.com/Deep-Agent/R1-V>. Accessed: 2025-02-02.

- Chen, Z.; Luo, X.; and Li, D. 2025. VisRL: Intention-Driven Visual Perception via Reinforced Reasoning. *arXiv Preprint*.
- Cheng, K.; Sun, Q.; Chu, Y.; Xu, F.; Li, Y.; Zhang, J.; and Wu, Z. 2024. SeeClick: Harnessing GUI Grounding for Advanced Visual GUI Agents. In *ACL*, volume 1, 9313–9332.
- DeepSeek-AI; Guo, D.; Yang, D.; Zhang, H.; Song, J.; and et al, R. Z. 2025. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. *arXiv Preprint*.
- Gou, B.; Wang, R.; Zheng, B.; Xie, Y.; Chang, C.; Shu, Y.; Sun, H.; and Su, Y. 2025. Navigating the Digital World as Humans Do: Universal Visual Grounding for GUI Agents. In *ICLR*.
- Hong, W.; Wang, W.; Lv, Q.; Xu, J.; Yu, W.; Ji, J.; Wang, Y.; Wang, Z.; Dong, Y.; Ding, M.; and Tang, J. 2024. Cog-Agent: A Visual Language Model for GUI Agents. In *CVPR*, 14281–14290.
- Jaech, A.; Kalai, A.; Lerer, A.; Richardson, A.; and et al, A. E. 2024. OpenAI o1 System Card. *arXiv Preprint*.
- Kwon, W.; Li, Z.; Zhuang, S.; Sheng, Y.; Zheng, L.; Yu, C. H.; Gonzalez, J.; Zhang, H.; and Stoica, I. 2023. Efficient Memory Management for Large Language Model Serving with PagedAttention. In *SOSP*, 611–626.
- Li, J.; and Huang, K. 2025. A Summary on GUI Agents with Foundation Models Enhanced by Reinforcement Learning. *arXiv Preprint*.
- Li, K.; Meng, Z.; Lin, H.; Luo, Z.; Tian, Y.; Ma, J.; Huang, Z.; and Chua, T. 2025. ScreenSpot-Pro: GUI Grounding for Professional High-Resolution Computer Use. *arXiv Preprint*.
- Li, W.; Bishop, W. E.; Li, A.; Rawles, C.; Campbell-Ajala, F.; Tyamagundlu, D.; and Riva, O. 2024. On the Effects of Data Scale on UI Control Agents. In *NeurIPS*.
- Lin, K. Q.; Li, L.; Gao, D.; Yang, Z.; Wu, S.; Bai, Z.; Lei, S. W.; Wang, L.; and Shou, M. Z. 2025. ShowUI: One Vision-Language-Action Model for GUI Visual Agent. In *CVPR*, 19498–19508.
- Liu, Y.; Li, P.; Xie, C.; Hu, X.; Han, X.; Zhang, S.; Yang, H.; and Wu, F. 2025a. InfiGUI-R1: Advancing Multimodal GUI Agents from Reactive Actors to Deliberative Reasoners. *arXiv Preprint*.
- Liu, Z.; Sun, Z.; Zang, Y.; Dong, X.; Cao, Y.; Duan, H.; Lin, D.; and Wang, J. 2025b. Visual-RFT: Visual Reinforcement Fine-Tuning. *arXiv Preprint*.
- Lu, Z.; Chai, Y.; Guo, Y.; Yin, X.; Liu, L.; Wang, H.; Xiong, G.; and Li, H. 2025. UI-R1: Enhancing Action Prediction of GUI Agents by Reinforcement Learning. *arXiv Preprint*.
- Luo, R.; Wang, L.; He, W.; and Xia, X. 2025. GUI-R1 : A Generalist R1-Style Vision-Language Action Model For GUI Agents. *arXiv Preprint*.
- Meng, F.; Du, L.; Liu, Z.; Zhou, Z.; Lu, Q.; Fu, D.; Shi, B.; Wang, W.; He, J.; Zhang, K.; Luo, P.; Qiao, Y.; Zhang, Q.; and Shao, W. 2025. MM-Eureka: Exploring Visual Aha Moment with Rule-based Large-scale Reinforcement Learning. *arXiv Preprint*.
- Peng, Y.; Zhang, G.; Zhang, M.; You, Z.; Liu, J.; Zhu, Q.; Yang, K.; Xu, X.; Geng, X.; and Yang, X. 2025. LMM-R1: Empowering 3B LMMs with Strong Reasoning Abilities Through Two-Stage Rule-Based RL. *arXiv Preprint*.
- Qin, Y.; Ye, Y.; Fang, J.; and et al, H. W. 2025. UI-TARS: Pioneering Automated GUI Interaction with Native Agents. *arXiv Preprint*.
- Rafailov, R.; Sharma, A.; Mitchell, E.; Manning, C. D.; Ermon, S.; and Finn, C. 2023. Direct Preference Optimization: Your Language Model is Secretly a Reward Model. In *NeurIPS*.
- Schulman, J.; Wolski, F.; Dhariwal, P.; Radford, A.; and Klimov, O. 2017. Proximal Policy Optimization Algorithms. *arXiv Preprint*.
- Shao, Z.; Wang, P.; Zhu, Q.; Xu, R.; Song, J.; Zhang, M.; Li, Y. K.; Wu, Y.; and Guo, D. 2024. DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models. *arXiv Preprint*.
- Shen, Y.; Zhang, J.; Huang, J.; Shi, S.; Zhang, W.; Yan, J.; Wang, N.; Wang, K.; and Lian, S. 2025. DAST: Difficulty-Adaptive Slow-Thinking for Large Reasoning Models. *arXiv Preprint*.
- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L.; and Polosukhin, I. 2017. Attention is All you Need. In *NIPS*, 5998–6008.
- von Werra, L.; Belkada, Y.; Tunstall, L.; Beeching, E.; Thrush, T.; Lambert, N.; Huang, S.; Rasul, K.; and Gallouédec, Q. 2020. TRL: Transformer Reinforcement Learning. <https://github.com/huggingface/trl>. Accessed: 2025-04-02.
- Wan, J.; Song, S.; Yu, W.; Liu, Y.; Cheng, W.; Huang, F.; Bai, X.; Yao, C.; and Yang, Z. 2024. OMNIPARSER: A Unified Framework for Text Spotting, Key Information Extraction and Table Recognition. In *CVPR*, 15641–15653.
- Wang, S.; Liu, W.; Chen, J.; Gan, W.; Zeng, X.; Yu, S.; Hao, X.; Shao, K.; Wang, Y.; and Tang, R. 2024. GUI Agents with Foundation Models: A Comprehensive Survey. *arXiv Preprint*.
- Wu, Z.; Wu, Z.; Xu, F.; Wang, Y.; Sun, Q.; Jia, C.; Cheng, K.; Ding, Z.; Chen, L.; Liang, P. P.; and Qiao, Y. 2025. OS-ATLAS: Foundation Action Model for Generalist GUI Agents. In *ICLR*.
- Xie, T.; Deng, J.; Li, X.; Yang, J.; and et al, H. W. 2025. Scaling Computer-Use Grounding via User Interface Decomposition and Synthesis. *arXiv Preprint*.
- Xu, Y.; Wang, Z.; Wang, J.; Lu, D.; Xie, T.; Saha, A.; Sahoo, D.; Yu, T.; and Xiong, C. 2024. Aguvis: Unified Pure Vision Agents for Autonomous GUI Interaction. *arXiv Preprint*.
- Yang, Y.; Li, D.; Yang, Y.; and et al, Z. L. 2025. GRPO for GUI Grounding Done Right. <https://huggingface.co/blog>HelloKKMe/grounding-r1>. Accessed: 2025-06-13.
- Yang, Y.; Wang, Y.; Li, D.; Luo, Z.; Chen, B.; Huang, C.; and Li, J. 2024. Aria-UI: Visual Grounding for GUI Instructions. *arXiv Preprint*.
- Zhang, C.; He, S.; Qian, J.; Li, B.; Li, L.; Qin, S.; Kang, Y.; Ma, M.; Liu, G.; Lin, Q.; Rajmohan, S.; Zhang, D.;

and Zhang, Q. 2025. Large Language Model-Brained GUI Agents: A Survey. *Trans. Mach. Learn. Res.*, 2025.

Zhou, Y.; Dai, S.; Wang, S.; Zhou, K.; Jia, Q.; and Xu, J. 2025. GUI-G1: Understanding R1-Zero-Like Training for Visual Grounding in GUI Agents. *arXiv Preprint*.