

# Generalist Scanner Meets Specialist Locator: A Synergistic Coarse-to-Fine Framework for Robust GUI Grounding

Anonymous ACL submission

## Abstract

Grounding natural language queries in graphical user interfaces (GUIs) presents a challenging task that requires models to comprehend diverse UI elements across various applications and systems, while also accurately predicting the spatial coordinates for the intended operation. To tackle this problem, we propose *GMS: Generalist Scanner Meets Specialist Locator*, a synergistic coarse-to-fine framework that effectively improves GUI grounding performance. *GMS* leverages the complementary strengths of general vision-language models (VLMs) and small, task-specific GUI grounding models by assigning them distinct roles within the framework. Specifically, the general VLM acts as a ‘Scanner’ to identify potential regions of interest, while the fine-tuned grounding model serves as a ‘Locator’ that outputs precise coordinates within these regions. This design is inspired by how humans perform GUI grounding, where the eyes scan the interface and the brain focuses on interpretation and localization. Our whole framework consists of five stages and incorporates hierarchical search with cross-modal communication to achieve promising prediction results. Experimental results on the ScreenSpot-Pro dataset show that while the ‘Scanner’ and ‘Locator’ models achieve only 2.0% and 3.7% accuracy respectively when used independently, their integration within *GMS* framework yields an overall accuracy of 35.7%, representing a 10× improvement. Additionally, *GMS* significantly outperforms other strong baselines under various settings, demonstrating its robustness and potential for general-purpose GUI grounding.

## 1 Introduction

Grounding natural language queries in graphical user interfaces (GUIs) requires models to predict accurate coordinates for user-specified actions, enabling applications in agent control, device automation, and accessibility (Wang et al., 2025a;

Nguyen et al., 2025; Zhang et al., 2025a; Tang et al., 2025b). As vision-language models (VLMs) advance in multimodal reasoning, GUI grounding emerges as a key benchmark for evaluating their interactive capabilities (Hui et al., 2025; Wang et al., 2025b; Li et al., 2025; Cheng et al., 2024; Liu et al., 2024). GUI grounding is challenging due to the diverse structures, styles, and semantics of interfaces across platforms. It requires fine-grained understanding of both textual and non-textual elements, dense visual layouts, and context-dependent functions, making accurate interpretation difficult (Li et al., 2025; Wu and Xie, 2024; Wu et al., 2025a).

Existing approaches can be broadly categorized into two groups: (i) Training-based methods either fine-tune base vision-language models, such as Qwen2-VL-7B, to directly predict grounding coordinates, or employ reinforcement learning techniques, such as GRPO, to induce multi-step reasoning processes that ultimately localize the target region (Gou et al., 2025; Wu et al., 2024; Shao et al., 2024). Although fine-tuning improves task-specific performance, it often sacrifices the model’s capacity for self-correction and adaptive reasoning. Reinforcement learning methods offer greater flexibility and generalization, but they incur substantial computational overhead and suffer from slow inference due to the complexity of the reasoning procedures they require (Luo et al., 2025a; Zhou et al., 2025; Lu et al., 2025; Liu et al., 2025; Tang et al., 2025a). (ii) Training-free methods seek to bypass the cost of retraining by leveraging pre-trained models. These include recursive zoom-in techniques that iteratively refine grounding predictions and planner-based strategies that utilize general models to guide localizers (Li et al., 2025; Wu et al., 2025a; Nguyen, 2025; Luo et al., 2025b; Ge et al., 2025; Zhang et al., 2024; Wang et al., 2024a). However, zoom-in strategies are highly sensitive to initial prediction errors and lack any verification mechanism, making them fragile in practice.

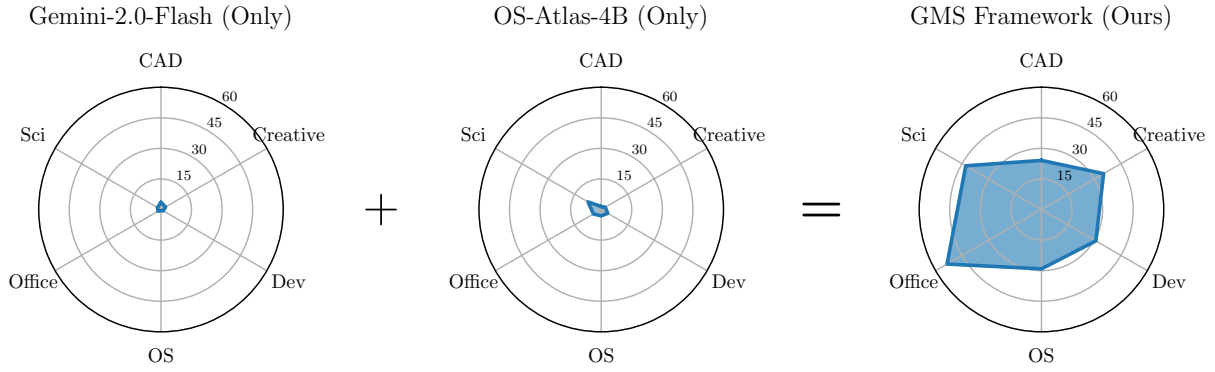


Figure 1: The two original models individually perform poorly on the GUI grounding task, with average accuracies below 4%. Under our GMS framework, where each model specializes in its strengths and collaborates effectively, the overall accuracy reaches 36%, which is nearly 10 $\times$  higher than their standalone performance.

Planner-based approaches mitigate this to some extent by introducing coordination between models, but they continue to rely on general models to produce bounding boxes. Since these general models are not trained explicitly for precise localization, the resulting predictions are often inaccurate, and errors tend to propagate throughout the grounding process.

To address these limitations, we propose **GMS: Generalist Scanner Meets Specialist Locator**, a synergistic coarse-to-fine framework. GMS integrates the complementary strengths of general-purpose and task-specific models to construct a training-free, modular grounding pipeline, as shown in Figure 2. The design of GMS is inspired by the human visual cognition process, in which broad perceptual scanning is followed by focused attention for fine-grained decision making. Accordingly, the general-purpose vision-language model operates as a ‘Scanner’ that identifies high-confidence candidate regions at a coarse level, while a fine-tuned GUI grounding model functions as a specialist ‘Locator’ that predicts precise coordinates within the selected regions. The GMS framework consists of five modules that enable coarse-to-fine localization: (1) *Hierarchical attention allocation*, where the ‘Scanner’ partitions the screen into coarse grids and selects semantically relevant regions; (2) *Iterative focus refinement*, where ambiguous areas are recursively zoomed in through semantically guided subdivision; (3) *Cross-modal verification*, where the ‘Locator’ proposes coordinates that are validated by the ‘Scanner’ to suppress false positives; (4) *Multi-agent consensus*, where the ‘Scanner’ and ‘Locator’ predictions are fused with asymmetric weighting for robust agreement; and (5) *Adaptive resolution enhancement*, where multi-scale late fusion rec-

onciles coarse semantic cues with fine pixel-level localization. This design creates a cognitively inspired perception and action loop, outperforming prior pipelines in both robustness and precision. In summary, our contributions are threefold:

(1) We introduce *GMS*, a training-free multi-agent framework that emulates human-like grounding by assigning complementary roles to generalist and specialist models, achieving substantial gains without additional fine-tuning.

(2) Experiments on the ScreenSpot-Pro dataset show that *GMS* improves performance by more than 10 $\times$  with weak model pairs and consistently outperforms other strong baselines, demonstrating both robustness and generalizability.

(3) We conduct extensive evaluations, including test-time scaling and ablations, to validate the framework. The results show that agents are most effective when specialized, leading to robust performance across diverse and challenging scenarios.

## 2 Related Works

Recent years witness substantial progress in GUI agent research, which evolves from rule-based web automation toward general-purpose control of user interfaces across diverse platforms, including mobile and desktop systems. A persistent challenge is the reliable localization of interface elements, which remains a key bottleneck for robust automation (Nakano et al., 2022; Zhang et al., 2025b; Wang et al., 2024a). The emergence of vision-language models marks a shift toward perception-driven grounding by leveraging both visual and textual inputs, without depending solely on structured metadata. Recent work improves robustness by fine-tuning VLMs on GUI-specific datasets, resulting in models that predict element coordinates with

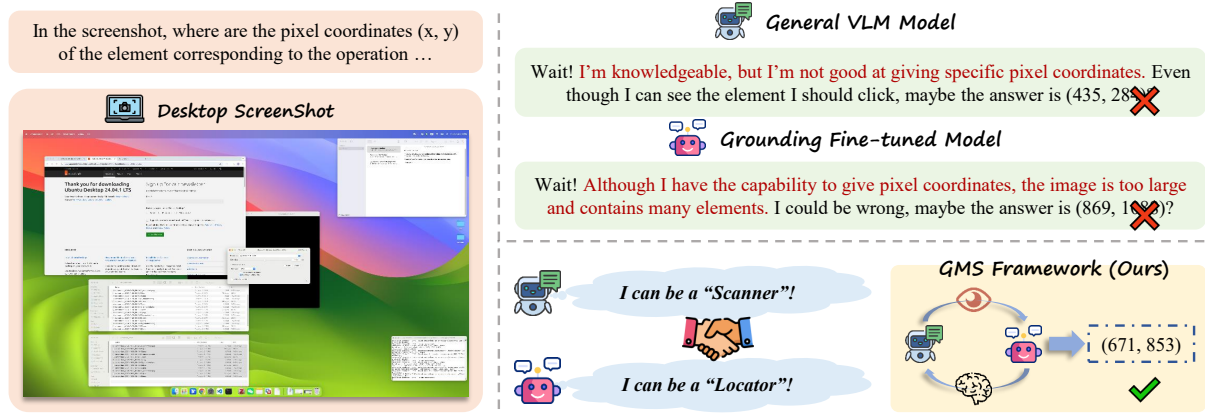


Figure 2: A simplified illustration conveys the motivation and high-level design of the proposed GMS framework.

higher precision (Gou et al., 2025; Wu et al., 2024; Gu et al., 2025; Qin et al., 2025). Some studies extend this direction using reinforcement learning approaches for multi-step decision-making with interpretable intermediate outputs (Tang et al., 2025c; Luo et al., 2025a; Wu et al., 2025b). In parallel, training-free approaches explore dual-system models, iterative zoom-in mechanisms, and the repurposing of general purpose models as planners to guide action selection (Wu et al., 2025a; Li et al., 2025). However, existing methods often overlook collaborative agent architectures, in which two specialized models assume distinct roles aligned with their respective strengths. Such cooperation presents a promising direction for integrating complementary model capabilities in GUI grounding.

### 3 Methodology

GUI grounding poses a dual challenge: it requires both global semantic understanding and precise spatial localization. Prior approaches often rely on a single model to handle both tasks simultaneously, leading to trade-offs that limit overall performance. Inspired by the dual-stream hypothesis in visual cognition, which separates the ‘what/where’ pathway from the ‘how’ pathway in human perception, we propose **GMS: Generalist Scanner Meets Specialist Locator**, a framework that explicitly decomposes the grounding task into two specialized agents: a generalist vision-language model acting as a ‘Scanner’, and a fine-tuned GUI grounding model serving as a ‘Locator’. GMS follows a coarse-to-fine strategy across five stages, with the detailed process illustrated in Figure 3.

Formally, let the GUI screen be an image  $I \in \mathbb{R}^{H \times W \times 3}$  and a natural language instruction as  $Q$ . The goal is to predict a pixel coordinate

$p = (x^*, y^*) \in [0, W] \times [0, H]$  corresponding to the GUI element described in  $Q$ . GMS achieves this through the following stages:

#### 3.1 Hierarchical Attention Allocation

Human visual attention operates in a coarse-to-fine manner, allocating cognitive resources hierarchically across the scene. Psychological studies show that within the first 300ms of exposure, humans can perform scene parsing to identify regions of interest prior to fine-scale analysis. GMS emulates this behavior via adaptive grid partitioning and region-level semantic scoring.

Specifically, we begin by decomposing the screen into a  $3 \times 3$  grid:

$$R = \{R_1, R_2, \dots, R_9\}, \quad R_i \subset I.$$

Each region  $R_i$  is defined by its bounding box  $B_i = [x_1^i, y_1^i, x_2^i, y_2^i]$ . The generalist vision-language model (e.g., GPT, Gemini) then acts as the ‘Scanner’, which evaluates each region’s semantic relevance to the query  $Q$  by computing:

$$s_i = \text{Select}(Inst_{\text{selection}}, Q, R_i), \quad s_i \in [0, 100].$$

The top- $k$  scoring regions are selected to form the candidate set  $R_{\text{top}}$ .

#### 3.2 Iterative Focus Refinement

One-shot attention allocation often fails in high-density GUI scenes due to:

- Semantic ambiguity from visually similar but functionally distinct elements.
- Contextual dependencies requiring reasoning over inter-element relations (e.g., ‘checkbox next to the password field’).

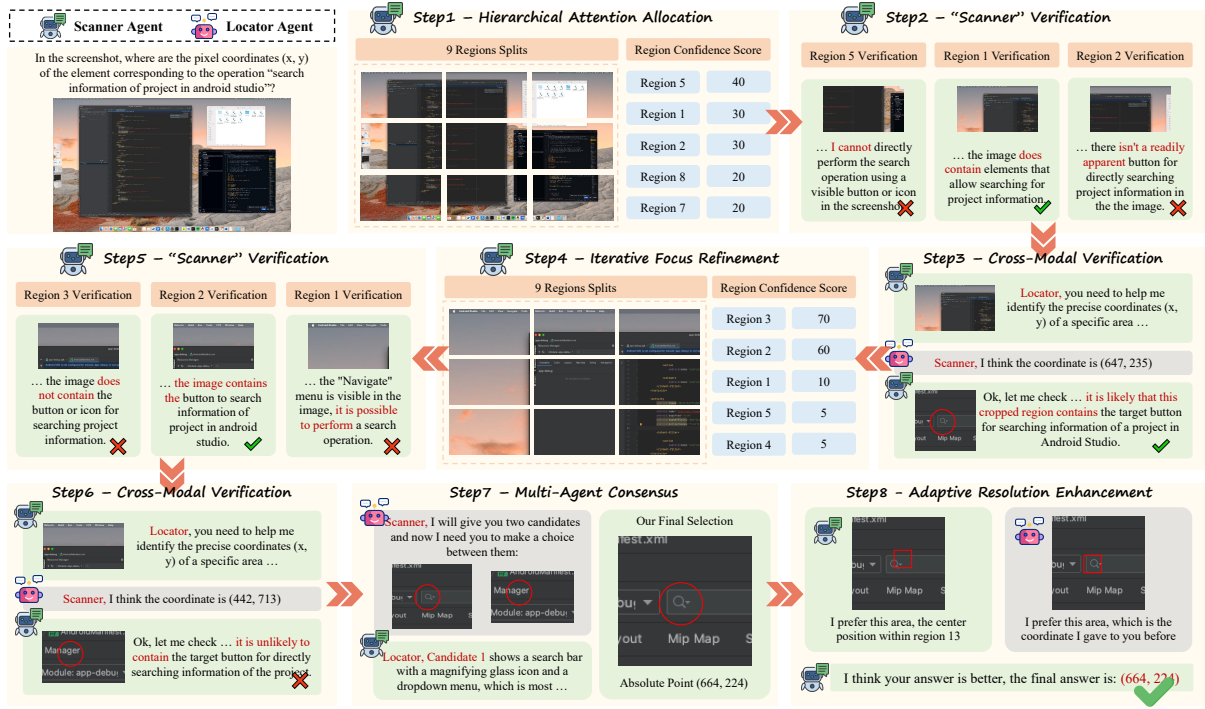


Figure 3: A detailed illustration of the proposed GMS framework highlights its multi-stage and hierarchical process. The ‘Scanner’ module mimics human vision by constraining the search space and identifying regions of interest, while the ‘Locator’ module emulates cognitive decision-making to determine precise coordinates.

To mitigate this, we design a recursive depth-first search (DFS) refinement process. At each level  $l$ , regions  $R^{(l)}$  are recursively subdivided into  $3 \times 3$  subgrids. The ‘Scanner’ re-applies the selection function:

$$R^{(l+1)} = \text{Select}(R^{(l)}, \text{Inst}_{\text{selection}}, Q),$$

until one of two stopping conditions is met: (i) the region’s width or height is below a threshold (e.g.,  $< 600$  px), or (ii) subsequent verification (section 3.3) indicates insufficient confidence.

### 3.3 Cross-Modal Verification

While generalist models excel at region-level semantic matching, they often suffer from false confidence due to hallucinations or overgeneralization. To correct this, we introduce a cross-modal verification mechanism that uses the specialist ‘Locator’ as a factuality check.

For each selected region  $R^{(l)}$ , the ‘Locator’ agent predicts a coordinate:

$$\hat{p}_l = \text{GUIGround}(Q, R^{(l)}).$$

A crop  $C_l$  of size  $125 \times 125$  pixels is extracted around  $\hat{p}_l$ , providing localized context. The ‘Scanner’ agent then performs verification:

$$v_l = \text{Verify}(C_l, Q, \text{Inst}_{\text{verification}}), \quad v_l \in \{0, 1\}.$$

Empirically,  $125 \times 125$  provides sufficient local cues while avoiding dilution from too many unrelated UI elements.

### 3.4 Multi-Agent Consensus

After multiple rounds of verification, we could obtain  $t$  candidate crops:

$$\mathcal{C} = \{(C_1, v_1), \dots, (C_t, v_t)\}, \quad v_l \in \{0, 1\}.$$

Selecting the best candidate is framed as a multi-agent consensus problem. Instead of naïve majority voting, we adopt an asymmetric weighting strategy, reflecting each agent’s relative expertise:

(i) The ‘Scanner’ agent contributes global context understanding and high-level semantic reasoning across multiple candidate regions.

(ii) The ‘Locator’ agent contributes fine-grained spatial precision and reliable confidence estimation within localized regions.

The ‘Scanner’ agent is instructed as follows:

$$\hat{l} \leftarrow \text{Eval}(\text{Inst}_{\text{evaluation}}, \mathcal{C}), \quad C^* = C_{\hat{l}}.$$

This step ensures that the selected region maximally aligns with both semantic and spatial constraints of the instruction  $Q$ .

Base Model	Development			Creative			CAD			Scientific			Office			OS			Average		
	text	icon	avg	text	icon	avg	text	icon	avg	text	icon	avg	text	icon	avg	text	icon	avg	text	icon	avg
GPT-4o	1.3	0.0	0.7	1.0	0.0	0.6	2.0	0.0	1.5	2.1	0.0	1.2	1.1	0.0	0.6	0.0	0.0	0.0	1.3	0.0	0.8
Gemini-2.0-Flash	0.6	2.1	1.3	4.5	0.0	2.6	3.6	3.1	3.4	2.1	1.8	2.0	2.3	0.0	1.7	0.0	1.1	0.5	2.5	1.3	2.0
Gemini-2.5-Flash-Lite	1.3	0.0	0.7	2.0	4.2	2.9	7.6	1.6	6.1	4.2	0.9	2.8	2.3	3.8	2.6	0.0	1.1	0.5	3.2	1.8	2.7
Claude (Computer Use)	22.0	3.9	12.6	25.9	3.4	16.8	14.5	3.7	11.9	33.9	15.8	25.8	30.1	16.3	26.2	11.0	4.5	8.1	23.4	7.1	17.1
UI-TARS-7B	58.4	12.4	36.1	50.0	9.1	32.8	20.8	9.4	18.0	63.9	31.8	50.0	63.3	20.8	53.5	30.8	16.9	24.5	47.8	16.2	35.7
UI-TARS-72B	63.0	17.3	40.8	57.1	15.4	39.6	18.8	12.5	17.2	64.6	20.9	45.7	63.3	26.4	54.8	42.1	15.7	30.1	50.9	17.5	38.1
OS-Atlas-4B	7.1	0.0	3.7	3.0	1.4	2.3	2.0	0.0	1.5	9.0	5.5	7.5	5.1	3.8	4.4	5.6	0.0	3.1	5.0	1.7	3.7
+ DiMo-GUI	13.6	1.4	7.7	9.6	2.8	6.7	4.1	4.7	4.2	30.6	4.5	19.3	24.3	15.1	22.2	7.5	2.2	5.1	14.6	4.0	10.6
+ GMS (w/ Gemini-2.0-Flash)	<b>44.2</b>	<u>16.6</u>	<b>30.8</b>	<b>49.0</b>	<b>16.1</b>	<b>35.2</b>	<b>27.9</b>	<b>12.5</b>	<b>24.1</b>	<b>56.3</b>	<b>25.5</b>	<b>42.9</b>	<b>57.6</b>	<b>39.6</b>	<b>53.5</b>	<b>36.4</b>	<b>20.2</b>	<b>29.1</b>	<b>45.2</b>	<b>20.2</b>	<b>35.7</b>
Δ	<b>37.1</b>	<b>16.6</b>	<b>27.1</b>	<b>46.0</b>	<b>14.7</b>	<b>32.9</b>	<b>25.9</b>	<b>12.5</b>	<b>22.6</b>	<b>47.3</b>	<b>20.0</b>	<b>35.4</b>	<b>52.5</b>	<b>35.8</b>	<b>49.1</b>	<b>30.8</b>	<b>20.2</b>	<b>26.0</b>	<b>40.2</b>	<b>18.5</b>	<b>32.0</b>
+ GMS (w/ Gemini-2.5-Flash-Lite)	<u>39.0</u>	<b>18.6</b>	<u>29.1</u>	<u>48.5</u>	<u>14.7</u>	<u>34.3</u>	<u>21.3</u>	<b>12.5</b>	<u>19.2</u>	<u>45.8</u>	<u>20.0</u>	<u>34.6</u>	<u>55.4</u>	<u>24.5</u>	<u>48.3</u>	<u>35.5</u>	<u>19.1</u>	<u>28.1</u>	<u>40.9</u>	<u>17.9</u>	<u>32.1</u>
Δ	<b>31.9</b>	<b>18.6</b>	<b>25.4</b>	<b>45.5</b>	<b>13.3</b>	<b>32.1</b>	<b>19.3</b>	<b>12.5</b>	<b>17.7</b>	<b>36.8</b>	<b>14.5</b>	<b>27.1</b>	<b>50.3</b>	<b>20.7</b>	<b>43.9</b>	<b>29.9</b>	<b>19.1</b>	<b>25.0</b>	<b>35.9</b>	<b>16.2</b>	<b>28.4</b>
UGround-7B	26.6	2.1	14.7	27.3	2.8	17.0	14.2	1.6	11.1	31.9	2.7	19.3	31.6	11.3	27.9	17.8	0.0	9.7	25.0	2.8	16.5
+ DiMo-GUI	44.2	6.2	25.8	39.9	7.7	26.4	17.3	3.1	13.8	50.7	8.2	32.3	46.9	15.1	39.6	32.7	10.1	22.4	38.1	7.9	26.6
+ GMS (w/ Qwen2.5-VL-7B)	44.2	13.8	29.4	56.1	15.4	39.0	<u>33.5</u>	<u>17.2</u>	<u>29.5</u>	<u>54.2</u>	25.5	41.7	59.3	34.0	53.5	37.4	18.0	28.6	47.9	19.0	36.9
Δ	<b>15.6</b>	<b>11.7</b>	<b>14.7</b>	<b>28.8</b>	<b>12.6</b>	<b>22.0</b>	<b>19.3</b>	<b>15.6</b>	<b>18.4</b>	<b>22.3</b>	<b>22.8</b>	<b>22.4</b>	<b>27.7</b>	<b>22.7</b>	<b>25.6</b>	<b>19.6</b>	<b>18.0</b>	<b>18.9</b>	<b>22.9</b>	<b>16.2</b>	<b>20.4</b>
+ GMS (w/ Gemini-2.0-Flash)	<b>60.4</b>	<b>24.1</b>	<b>42.8</b>	<b>63.1</b>	<b>24.5</b>	<b>46.9</b>	<b>35.5</b>	14.1	<b>30.3</b>	<b>62.5</b>	<u>27.3</u>	47.2	<b>71.8</b>	<b>43.4</b>	<b>65.2</b>	<b>52.3</b>	<b>27.0</b>	<b>40.8</b>	<b>57.4</b>	<b>25.8</b>	<b>45.4</b>
Δ	<b>33.8</b>	<b>22.0</b>	<b>28.1</b>	<b>35.8</b>	<b>21.7</b>	<b>29.9</b>	<b>21.3</b>	<b>12.5</b>	<b>19.2</b>	<b>30.6</b>	<b>24.6</b>	<b>27.9</b>	<b>40.2</b>	<b>32.1</b>	<b>37.3</b>	<b>34.5</b>	<b>27.0</b>	<b>31.1</b>	<b>32.4</b>	<b>23.0</b>	<b>28.9</b>
+ GMS (w/ Gemini-2.5-Flash-Lite)	<u>44.8</u>	<u>18.6</u>	<u>32.1</u>	<u>59.6</u>	<u>21.7</u>	<u>43.7</u>	29.9	<b>18.8</b>	27.2	50.0	<b>28.2</b>	40.6	<u>70.1</u>	<u>35.8</u>	<u>62.2</u>	<u>44.9</u>	<u>20.2</u>	<u>33.7</u>	<u>50.2</u>	<u>22.8</u>	<u>39.7</u>
Δ	<b>18.2</b>	<b>16.5</b>	<b>17.4</b>	<b>32.3</b>	<b>18.9</b>	<b>26.7</b>	<b>15.7</b>	<b>17.2</b>	<b>16.1</b>	<b>18.1</b>	<b>25.5</b>	<b>21.3</b>	<b>38.5</b>	<b>24.5</b>	<b>34.3</b>	<b>27.1</b>	<b>20.2</b>	<b>24.0</b>	<b>25.2</b>	<b>20.0</b>	<b>23.2</b>
UGround-V1-7B	51.9	3.4	28.4	48.0	9.1	31.7	20.0	1.6	15.3	57.6	16.4	39.8	61.6	13.2	50.4	37.4	7.9	25.0	45.6	8.4	31.4
+ DiMo-GUI	57.8	21.4	40.1	60.1	18.1	42.5	45.7	18.8	39.1	<b>75.7</b>	28.2	55.1	<b>79.7</b>	37.7	70.0	<u>51.4</u>	<b>30.3</b>	<u>41.8</u>	61.7	24.3	47.4
+ GMS (w/ Qwen2.5-VL-7B)	53.2	20.7	37.5	57.1	19.6	41.3	59.4	<b>29.7</b>	52.1	62.5	34.5	50.4	67.8	35.8	60.4	45.8	15.7	32.1	58.4	<u>24.5</u>	45.5
Δ	<b>1.3</b>	<b>17.3</b>	<b>9.1</b>	<b>9.1</b>	<b>10.5</b>	<b>9.6</b>	<b>39.4</b>	<b>28.1</b>	<b>36.8</b>	<b>4.9</b>	<b>18.1</b>	<b>10.6</b>	<b>6.2</b>	<b>22.6</b>	<b>10.0</b>	<b>8.4</b>	<b>7.8</b>	<b>7.1</b>	<b>12.8</b>	<b>16.1</b>	<b>14.1</b>
+ GMS (w/ Gemini-2.0-Flash)	<b>69.5</b>	<b>35.9</b>	<b>53.2</b>	<b>67.7</b>	<b>26.6</b>	<b>50.4</b>	<b>67.0</b>	<u>28.1</u>	<b>57.5</b>	<u>70.1</u>	<b>38.2</b>	<b>56.3</b>	76.8	<b>50.9</b>	<b>70.9</b>	<u>51.4</u>	21.3	37.8	68.1	32.5	<u>54.5</u>
Δ	<b>17.6</b>	<b>32.5</b>	<b>24.8</b>	<b>19.7</b>	<b>17.5</b>	<b>18.7</b>	<b>47.0</b>	<b>26.5</b>	<b>42.2</b>	<b>12.5</b>	<b>21.8</b>	<b>16.5</b>	<b>15.2</b>	<b>37.7</b>	<b>20.5</b>	<b>14.0</b>	<b>13.4</b>	<b>12.8</b>	<b>22.5</b>	<b>24.1</b>	<b>23.1</b>
+ GMS (w/ Gemini-2.5-Flash-Lite)	<u>59.1</u>	<u>29.7</u>	44.8	<b>72.2</b>	<b>30.8</b>	<b>54.8</b>	<b>70.6</b>	17.2	<b>57.5</b>	69.4	<b>38.2</b>	<u>55.9</u>	<u>78.0</u>	<u>45.3</u>	<u>70.4</u>	<b>57.9</b>	<u>29.2</u>	<b>44.9</b>	<b>68.9</b>	<b>31.5</b>	<b>54.6</b>
Δ	<b>7.2</b>	<b>26.3</b>	<b>16.4</b>	<b>24.2</b>	<b>21.7</b>	<b>23.1</b>	<b>50.6</b>	<b>15.6</b>	<b>42.2</b>	<b>11.8</b>	<b>21.8</b>	<b>16.1</b>	<b>16.4</b>	<b>32.1</b>	<b>20.0</b>	<b>20.5</b>	<b>21.3</b>	<b>19.9</b>	<b>23.3</b>	<b>23.1</b>	<b>23.2</b>

Table 1: Main experimental results on the ScreenSpot-Pro dataset. The table reports performance under the proposed GMS framework with different combinations of ‘Scanner’ and ‘Locator’ agents, compared against a range of baseline methods. The best accuracy for each setting is highlighted in **bold**, and the second-best is underlined. Relative improvements (in percentage points) are annotated.

### 3.5 Adaptive Resolution Enhancement

The final prediction requires resolving discrepancies between coarse attention and fine spatial cues. Generalist vision-language models operate on low-resolution patches, while the specialist operates on raw pixels. To bridge this, we design a multi-scale late fusion module.

First, we upscale  $C^*$  by  $\times 5$  in both dimensions to improve resolution. A  $5 \times 5$  grid is imposed, followed by a  $3 \times 3$  subgrid within the selected region. The ‘Scanner’ estimates a coarse point:

$$p_{\text{scanner}} = \text{Center}(z^*).$$

In parallel, the ‘Locator’ provides a direct prediction:

$$p_{\text{locator}} = \text{GUIGround}(Q, C^*).$$

The final decision is delegated to the stronger ‘Scanner’ agent:

$$p_{\text{final}} = \text{Decide}(Q, C^*, \{p_{\text{scanner}}, p_{\text{locator}}\}, \text{Inst}_{\text{decision}})$$

This fusion mechanism leverages multi-resolution reasoning, ensuring that the final coordinate prediction is both semantically coherent and spatially precise.

## 4 Experiments Setup

### 4.1 Dataset

We evaluate our framework on the **ScreenSpot-Pro** benchmark, which consists of over 1,500 high-resolution desktop screenshots spanning six GUI grounding tasks (Li et al., 2025).

### 4.2 Vision Language Models

We instantiate our framework with two vision-language models in complementary roles: a general-purpose ‘Scanner’ for broad visual understanding and instruction following, and a GUI-specialized ‘Locator’ for precise element localization.

To demonstrate that the framework effectively exploits each model’s strengths, we select two well-known grounding-focused model families, each

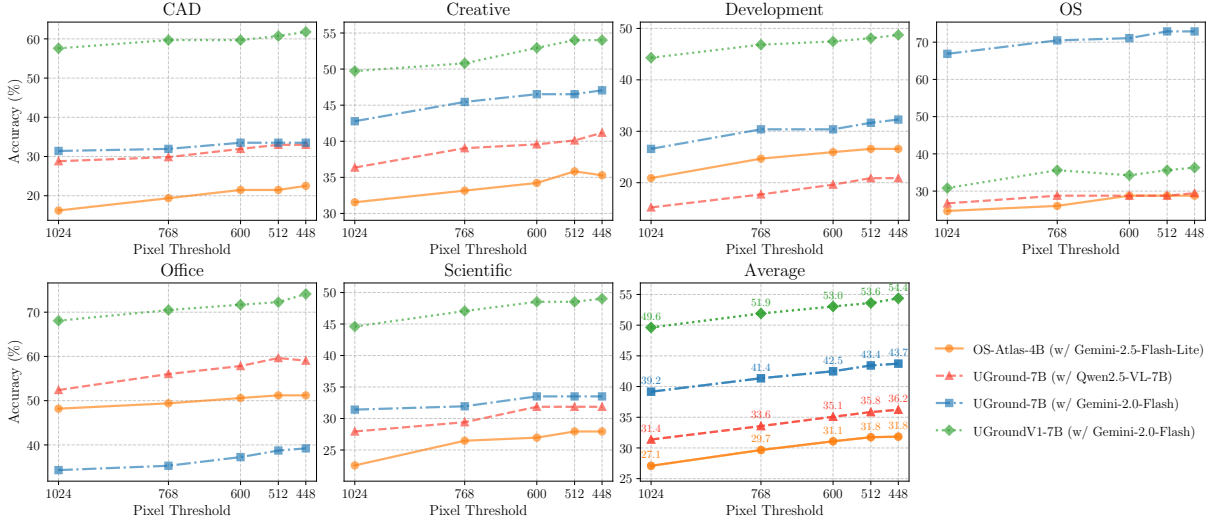


Figure 4: Experimental results illustrating the impact of decreasing the pixel number threshold from 1024 to 448 in the hierarchical search constraints. The figure reports the accuracy of six sub-categories and the overall accuracy across four agentic combinations.

with fewer than 7B parameters: OS-Atlas (Wu et al., 2024) and UGround (Gou et al., 2025).

For the ‘Scanner’ role, we balance cost and model availability (including both open-weight and closed-source models) and choose: Qwen2.5-VL (Bai et al., 2025) and Gemini (Google, 2025; Comanici et al., 2025; Team et al., 2024, 2025b).

### 4.3 Metrics

We use **accuracy** as the evaluation metric. Formally, let  $\hat{\mathbf{p}} = (x, y)$  denote the predicted coordinate and  $\mathcal{B} = [x_{\min}, x_{\max}] \times [y_{\min}, y_{\max}]$  denote the ground-truth bounding box. Then, the accuracy over  $N$  samples is:  $Accuracy = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(\hat{\mathbf{p}}_i \in \mathcal{B}_i)$ .

### 4.4 Implementation Details

We obtain all open-weight models from their official repositories on HuggingFace. For these fine-tuned grounding models, we set the temperature to 0.0 to ensure faithfulness. For closed-source models, we perform inference via the OpenRouter platform. To ensure consistency and efficiency, we adopt the default inference settings: temperature = 0.7 and top<sub>p</sub> = 0.95. The prompts and baseline introduction are provided in Appendix G and Appendix E, respectively.

## 5 Experiment Results

We evaluate our proposed framework on the ScreenSpot-Pro benchmark, with results presented in Table 1. Our framework consistently outperforms all baselines across multiple settings, includ-

ing strong fine-tuned models (up to 72B parameters) and DiMo-GUI. The improvements are particularly substantial across various sub-categories, covering both text and icon grounding tasks, with relative gains ranging from 100% to over 1000%. We highlight the following key findings:

**Effectiveness in Low-Performance Settings.** The OS-Atlas-4B model performs poorly under direct inference, achieving only 13% accuracy even with DiMo-GUI. Remarkably, when integrated into our framework and paired with Gemini models, each of which yields less than 3% accuracy individually, the combined system achieves 30% accuracy. This represents a 10× improvement over the original results and a 2× improvement over DiMo-GUI. These findings highlight the framework’s ability to coordinate weaker models into a highly effective cooperative system by assigning them specialized roles.

**Superior Performance on Icon Grounding Tasks.** Existing methods, including DiMo-GUI and large fine-tuned models, typically underperform on icon-related grounding tasks compared to text-based tasks. In contrast, our GMS framework substantially alleviates this disparity. By leveraging the synergy between the generalist ‘Scanner’ and specialist ‘Locator’ modules, our method boosts icon grounding accuracy from 1.7% to 20% on OS-Atlas-4B (a 1076% improvement) and from 2.8% to 25.8% on UGround-7B (an 821% improvement).

**Scalability and Model Flexibility.** Our frame-

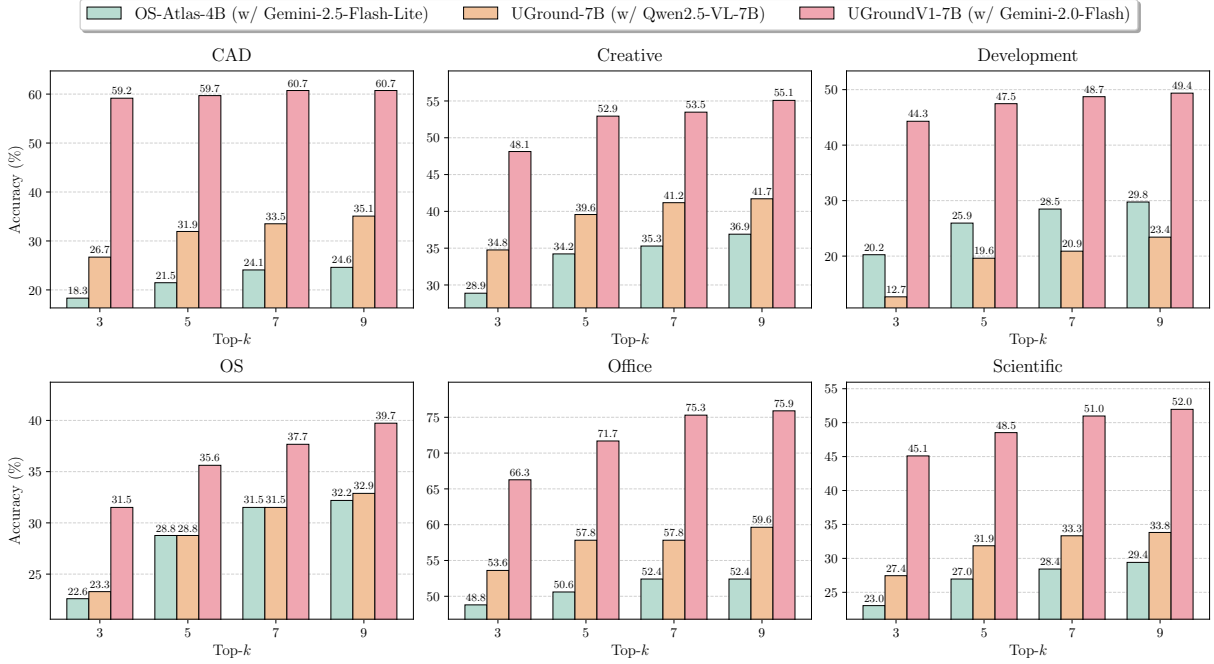


Figure 5: Experimental results showing the impact of increasing the top- $k$  selection value from 3 to 9. The figure reports the accuracy across six sub-categories under four agentic combinations.

work demonstrates strong scalability. Even with relatively small ‘Scanner’ models, such as 7B or flash-Gemini variants, the performance gains remain substantial. Furthermore, the results indicate that stronger general models lead to better outcomes. Given the framework’s flexible and modular design, integrating more capable models (e.g., stronger Gemini variants) may yield further performance improvements.

## 6 Test-Time Scaling

Our framework integrates the concept of test-time scaling into the hierarchical search process. The inference time can be flexibly controlled by adjusting two key factors: the top- $k$  selection parameter at each search step and the pixel threshold used during the process. In this section, we analyze the impact of these parameters on the 15 most challenging subsets.

### 6.1 Impact of Pixel Value Threshold

We begin by evaluating the impact of different pixel thresholds, with accuracy results shown in Figure 4. As the threshold decreases from 1024 to 512, which allows the ‘Scanner’ agents to capture finer-grained details of the GUI interface, we observe a consistent increase in accuracy. Notably, even at the higher threshold of 1024, the framework maintains strong performance, with only a marginal drop in

accuracy compared to lower thresholds. This result highlights the robustness and effectiveness of our proposed framework. Moreover, the improvements observed with decreasing thresholds suggest promising test-time scaling capabilities. Given that many original images in the ScreenSpot-Pro dataset exceed 3000 pixels in resolution, these findings indicate that the framework is well suited for high-resolution settings, which are critical in GUI understanding tasks.

### 6.2 Impact of Top- $k$ Region Selection

We further examine the impact of the  $k$ -value in the top- $k$  region selection mechanism, which guides the deeper stages of hierarchical search. As shown in Figure 5, increasing  $k$  from 3 to 9 leads to a gradual improvement in accuracy. This trend aligns with intuition, as a larger search space allows the ‘Scanner’ agent to explore more potentially relevant subregions. Since in high-resolution settings, where input images are especially large, the agent may fail to cover all important areas. A smaller  $k$  can lead to the omission of critical regions, particularly when the agent assigns low confidence to relevant areas due to limited context or weaker understanding.

The most significant performance gain occurs when increasing  $k$  from 3 to 5, suggesting that the ‘Scanner’ agent possesses a baseline level of capability, and a modest expansion of the search space

Scanner Agent	Locator Agent	Inference Method	Overall Accuracy (%)
Gemini-2.0-Flash	OS-Atlas-4B	GMS	35.67
		<i>w/o Cross-Modal Verification</i>	27.83(↓ 7.84)
		<i>w/o Multi-Agent Consensus</i>	33.05(↓ 2.62)
		<i>w/o Adaptive Resolution Enhancement</i>	31.59(↓ 4.08)
Gemini-2.5-Flash-Lite	UGround-7B	GMS	39.72
		<i>w/o Cross-Modal Verification</i>	33.57(↓ 6.15)
		<i>w/o Multi-Agent Consensus</i>	37.20(↓ 2.52)
		<i>w/o Adaptive Resolution Enhancement</i>	34.91(↓ 4.81)

Table 2: Ablation results after individually removing each of the three crucial stages from our framework. Overall accuracy drops significantly compared to the full framework, highlighting the cooperative nature and effectiveness of the proposed architecture.

greatly enhances its effectiveness. Beyond this point, further improvements are observed, but with diminishing returns. We also find that the choice of ‘Scanner’ agent influences the model’s sensitivity to changes in  $k$ . For example, the Qwen2.5-VL-7B model shows the most pronounced improvement as  $k$  increases; on the *Development* split, accuracy rises from 12.7% to 23.4%. In contrast, when stronger Gemini models are used as the ‘Scanner’ agent, the benefit of increasing  $k$  is less substantial. This aligns with expectations, as stronger vision-language models are generally more confident and accurate in identifying the correct region, thereby reducing the need for a large search space.

## 7 Ablation Study

To validate the effectiveness of each key component in our proposed framework, we conduct a series of ablation studies using three modified baselines:

- Deletion of Cross-Modal Verification: The ‘Locator’ agent’s predicted coordinates are passed directly, without any confidence-based filtering.
- Deletion of Multi-Agent Consensus: The ‘Locator’ agent always selects the coordinate with the highest confidence score, bypassing the consensus selection mechanism.
- Deletion of Adaptive Resolution Enhancement: The ‘Locator’ and ‘Scanner’ agents no longer collaborate; the output of the ‘Locator’ alone is used as the final click instruction.

The experimental results, presented in Table 2, show that removing any component leads to noticeable performance degradation. Among these experiments, the removal of cross-modal verification has the most severe impact. This is likely due

to the absence of confidence filtering, which causes the ‘Scanner’ agent to receive multiple candidate regions without sufficient guidance, making effective reasoning difficult, especially under long-context constraints or limited output reasoning capacity.

In contrast, the removal of multi-agent consensus has the least impact on performance. This finding reflects the strong evaluation capabilities of the ‘Scanner’ agent, which can reliably assess each candidate to determine whether it contains the target position. This further supports the design intuition of assigning distinct and complementary responsibilities to different agent types within our framework.

## 8 Conclusion

We propose GMS, a synergistic coarse-to-fine framework that employs hierarchical search with test-time scaling. Drawing inspiration from how humans approach GUI grounding tasks, GMS introduces two specialized roles: the ‘Scanner’ and the ‘Locator’. This division enables cooperative behavior, allowing each agent to focus on the subtask that aligns with its respective strengths. The ‘Scanner’ performs coarse region localization, while the ‘Locator’ is responsible for precise coordinate prediction within the identified region. The framework not only surpasses strong baselines but also substantially boosts the performance of two relatively weak models when combined, achieving nearly double their original accuracy without any additional fine-tuning. Ablation studies further confirm the robustness of the framework, showing that each stage contributes significantly to overall performance. These results underscore the practical applicability of GMS to real-world GUI grounding scenarios and highlight its potential as a general paradigm for agent collaboration in vision–language tasks.

## 498 Limitations

499 In this paper, we propose GMS, a synergistic  
500 coarse-to-fine framework that achieves significant  
501 performance improvements without requiring any  
502 additional fine-tuning. Despite the promising re-  
503 sults, several limitations remain: (1) Due to bud-  
504 get constraints and the large image sizes in the  
505 experimental dataset, we are currently unable to  
506 evaluate the framework on the full dataset using  
507 more diverse and powerful models, such as GPT-4o  
508 or Claude Opus (OpenAI et al., 2024; Anthropic,  
509 2024, 2025). (2) The inference prompts for both  
510 agents, particularly the ‘Scanner’, are manually  
511 crafted to ensure the reliability and relevance of the  
512 outputs. This introduces a dependency on prompt  
513 quality, where suboptimal prompt design can nega-  
514 tively affect performance. As part of future work,  
515 we plan to explore prompt-free agent frameworks,  
516 such as those based on attention mechanisms or  
517 other adaptive planning strategies.

## 518 Ethics Statement

519 Ethical considerations play a central role in this  
520 research. All models used in this study are either  
521 open-weight or widely adopted within the scienti-  
522 fic community, ensuring transparency and repro-  
523 ducibility. The proposed GMS framework aims  
524 to advance the capabilities of current VLM agents  
525 for the GUI grounding task, contributing to real-  
526 world applications without introducing or reinforc-  
527 ing harmful biases. No personally identifiable in-  
528 formation or sensitive data is involved in this work.  
529 We are committed to responsible research practices  
530 and advocate for the transparent reporting and ethi-  
531 cal deployment of AI technologies in ways that  
532 serve the broader interests of society.

## 533 References

534 Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien  
535 Bubeck, Ronen Eldan, Suriya Gunasekar, Michael  
536 Harrison, Russell J. Hewett, Mojan Javaheripi, Piero  
537 Kauffmann, James R. Lee, Yin Tat Lee, Yuanzhi Li,  
538 Weishung Liu, Caio C. T. Mendes, Anh Nguyen,  
539 Eric Price, Gustavo de Rosa, Olli Saarikivi, and  
540 8 others. 2024. *Phi-4 technical report*. *Preprint*,  
541 arXiv:2412.08905.

542 Anthropic. 2024. Claude 3.5 sonnet. [https://www.  
543 anthropic.com/news/claude-3-5-sonnet](https://www.anthropic.com/news/claude-3-5-sonnet). Ac-  
544 cessed Jun 20, 2024.

Anthropic. 2025. Introducing claude 4. [https://www.  
545 anthropic.com/news/claude-4](https://www.anthropic.com/news/claude-4). Accessed May  
546 22, 2025. 547

Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang,  
548 Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou,  
549 and Jingren Zhou. 2023. *Qwen-vl: A versatile vision-  
550 language model for understanding, localization, text  
551 reading, and beyond*. *Preprint*, arXiv:2308.12966. 552

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wen-  
553 bin Ge, Sibao Song, Kai Dang, Peng Wang, Shi-  
554 jie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu,  
555 Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei  
556 Wang, Wei Ding, Zheren Fu, Yiheng Xu, and 8 oth-  
557 ers. 2025. *Qwen2.5-vl technical report*. *Preprint*,  
558 arXiv:2502.13923. 559

Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu,  
560 Yantao Li, Jianbing Zhang, and Zhiyong Wu. 2024.  
561 *SeeClick: Harnessing gui grounding for advanced  
562 visual gui agents*. *Preprint*, arXiv:2401.10935. 563

Gheorghe Comanici, Eric Bieber, Mike Schaeckermann,  
564 Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Mar-  
565 cel Blistein, Ori Ram, Dan Zhang, Evan Rosen,  
566 Luke Marris, Sam Petulla, Colin Gaffney, Asaf Ahar-  
567 onni, Nathan Lintz, Tiago Cardal Pais, Henrik Ja-  
568 cobsson, Idan Szpektor, Nan-Jiang Jiang, and 2 oth-  
569 ers. 2025. *Gemini 2.5: Pushing the frontier with  
570 advanced reasoning, multimodality, long context,  
571 and next generation agentic capabilities*. *Preprint*,  
572 arXiv:2507.06261. 573

Haonan Ge, Yiwei Wang, Ming-Hsuan Yang, and Yu-  
574 jun Cai. 2025. *Mrfd: Multi-region fusion decoding  
575 with self-consistency for mitigating hallucinations in  
576 vlms*. *Preprint*, arXiv:2508.10264. 577

Google. 2025. Gemini 2.0. [https:  
578 //developers.googleblog.com/en/  
579 gemini-2-family-expands/](https://developers.googleblog.com/en/gemini-2-family-expands/). Accessed FEB.  
580 5, 2025. 581

Boyu Gou, Ruohan Wang, Boyuan Zheng, Yanan Xie,  
582 Cheng Chang, Yiheng Shu, Huan Sun, and Yu Su.  
583 2025. *Navigating the digital world as humans do:  
584 Universal visual grounding for gui agents*. *Preprint*,  
585 arXiv:2410.05243. 586

Zhangxuan Gu, Zhengwen Zeng, Zhenyu Xu, Xingran  
587 Zhou, Shuheng Shen, Yunfei Liu, Beitong Zhou,  
588 Changhua Meng, Tianyu Xia, Weizhi Chen, Yue Wen,  
589 Jingya Dou, Fei Tang, Jinzhen Lin, Yulin Liu, Zhen-  
590 lin Guo, Yichen Gong, Heng Jia, Changlong Gao,  
591 and 5 others. 2025. *Ui-venus technical report: Build-  
592 ing high-performance ui agents with rft*. *Preprint*,  
593 arXiv:2508.10833. 594

Wenyi Hong, Weihang Wang, Qingsong Lv, Jiazheng  
595 Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan  
596 Wang, Yuxuan Zhang, Juanzi Li, Bin Xu, Yuxiao  
597 Dong, Ming Ding, and Jie Tang. 2024. *Cogagent:  
598 A visual language model for gui agents*. *Preprint*,  
599 arXiv:2312.08914. 600



713	Da Yin, Hao Zeng, Jiajie Zhang, Kedong Wang,	Hang Wu, Hongkai Chen, Yujun Cai, Chang Liu, Qing-	772
714	Lucen Zhong, Mingdao Liu, Rui Lu, Shulin Cao,	wen Ye, Ming-Hsuan Yang, and Yiwei Wang. 2025a.	773
715	Xiaohan Zhang, Xuancheng Huang, Yao Wei, and	<a href="#">Dimo-gui: Advancing test-time scaling in gui ground-</a>	774
716	21 others. 2025a. <a href="#">Glm-4.5: Agentic, reasoning,</a>	<a href="#">ing via modality-aware visual reasoning</a> . <i>Preprint</i> ,	775
717	<a href="#">and coding (arc) foundation models</a> . <i>Preprint</i> ,	arXiv:2507.00008.	776
718	arXiv:2508.06471.		
719	Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-	Penghao Wu and Saining Xie. 2024. V?: Guided visual	777
720	Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan	search as a core mechanism in multimodal llms. In	778
721	Schalkwyk, Andrew M. Dai, Anja Hauth, Katie	<i>Proceedings of the IEEE/CVF Conference on Com-</i>	779
722	Millican, David Silver, Melvin Johnson, Ioannis	<i>puter Vision and Pattern Recognition (CVPR)</i> , pages	780
723	Antonoglou, Julian Schrittwieser, Amelia Glaese,	13084–13094.	781
724	Jilin Chen, Emily Pitler, Timothy Lillicrap, Angeliki		
725	Lazaridou, and 13 others. 2025b. <a href="#">Gemini: A fam-</a>	Qianhui Wu, Kanzhi Cheng, Rui Yang, Chaoyun Zhang,	782
726	<a href="#">ily of highly capable multimodal models</a> . <i>Preprint</i> ,	Jianwei Yang, Huiqiang Jiang, Jian Mu, Baolin	783
727	arXiv:2312.11805.	Peng, Bo Qiao, Reuben Tan, Si Qin, Lars Liden,	784
		Qingwei Lin, Huan Zhang, Tong Zhang, Jianbing	785
728	Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan	Zhang, Dongmei Zhang, and Jianfeng Gao. 2025b.	786
729	Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer,	<a href="#">Gui-actor: Coordinate-free visual grounding for gui</a>	787
730	Damien Vincent, Zhufeng Pan, Shibo Wang, Soroosh	<a href="#">agents</a> . <i>Preprint</i> , arXiv:2506.03143.	788
731	Mariooryad, Yifan Ding, Xinyang Geng, Fred Al-		
732	cobber, Roy Frostig, Mark Omernick, Lexi Walker,	Zhiyong Wu, Zhenyu Wu, Fangzhi Xu, Yian Wang,	789
733	Cosmin Paduraru, Christina Sorokin, and 12 oth-	Qishi Sun, Chengyou Jia, Kanzhi Cheng, Zichen	790
734	ers. 2024. <a href="#">Gemini 1.5: Unlocking multimodal un-</a>	Ding, Liheng Chen, Paul Pu Liang, and Yu Qiao.	791
735	<a href="#">derstanding across millions of tokens of context</a> .	2024. <a href="#">Os-atlas: A foundation action model for gener-</a>	792
736	<i>Preprint</i> , arXiv:2403.05530.	<a href="#">alist gui agents</a> . <i>Preprint</i> , arXiv:2410.23218.	793
737	Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya	Yuhao Yang, Yue Wang, Dongxu Li, Ziyang Luo, Bei	794
738	Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin,	Chen, Chao Huang, and Junnan Li. 2025. <a href="#">Aria-</a>	795
739	Tatiana Matejovicova, Alexandre Ramé, Morgane	<a href="#">ui: Visual grounding for gui instructions</a> . <i>Preprint</i> ,	796
740	Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey	arXiv:2412.16256.	797
741	Cideron, Jean bastien Grill, Sabela Ramos, Edouard		
742	Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev,	Chaoyun Zhang, Shilin He, Jiaxu Qian, Bowen Li,	798
743	and 14 others. 2025c. <a href="#">Gemma 3 technical report</a> .	Liqun Li, Si Qin, Yu Kang, Minghua Ma, Guyue Liu,	799
744	<i>Preprint</i> , arXiv:2503.19786.	Qingwei Lin, Saravan Rajmohan, Dongmei Zhang,	800
		and Qi Zhang. 2025a. <a href="#">Large language model-brained</a>	801
745	Junyang Wang, Haiyang Xu, Jiabo Ye, Ming Yan,	<a href="#">gui agents: A survey</a> . <i>Preprint</i> , arXiv:2411.18279.	802
746	Weizhou Shen, Ji Zhang, Fei Huang, and Jitao Sang.		
747	2024a. <a href="#">Mobile-agent: Autonomous multi-modal mo-</a>	Chaoyun Zhang, Liqun Li, Shilin He, Xu Zhang,	803
748	<a href="#">bile device agent with visual perception</a> . <i>Preprint</i> ,	Bo Qiao, Si Qin, Minghua Ma, Yu Kang, Qing-	804
749	arXiv:2401.16158.	wei Lin, Saravan Rajmohan, Dongmei Zhang, and	805
		Qi Zhang. 2024. <a href="#">Ufo: A ui-focused agent for win-</a>	806
750	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhi-	<a href="#">dows os interaction</a> . <i>Preprint</i> , arXiv:2402.07939.	807
751	hao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin		
752	Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei	Chi Zhang, Zhao Yang, Jiaxuan Liu, Yanda Li, Yucheng	808
753	Du, Xuancheng Ren, Rui Men, Dayiheng Liu,	Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu.	809
754	Chang Zhou, Jingren Zhou, and Junyang Lin. 2024b.	2025b. <a href="#">Appagent: Multimodal agents as smartphone</a>	810
755	<a href="#">Qwen2-vl: Enhancing vision-language model’s per-</a>	<a href="#">users</a> . In <i>Proceedings of the 2025 CHI Conference</i>	811
756	<a href="#">ception of the world at any resolution</a> . <i>Preprint</i> ,	<i>on Human Factors in Computing Systems</i> , CHI ’25,	812
757	arXiv:2409.12191.	New York, NY, USA. Association for Computing	813
		Machinery.	814
758	Shuai Wang, Weiwen Liu, Jingxuan Chen, Yuqi Zhou,	Qiyuan Zhang, Fuyuan Lyu, Zexu Sun, Lei Wang,	815
759	Weinan Gan, Xingshan Zeng, Yuhan Che, Shuai	Weixu Zhang, Wenyue Hua, Haolun Wu, Zhihan Guo,	816
760	Yu, Xinlong Hao, Kun Shao, Bin Wang, Chuhan	Yufei Wang, Niklas Muennighoff, Irwin King, Xue	817
761	Wu, Yasheng Wang, Ruiming Tang, and Jianye Hao.	Liu, and Chen Ma. 2025c. <a href="#">A survey on test-time</a>	818
762	2025a. <a href="#">Gui agents with foundation models: A com-</a>	<a href="#">scaling in large language models: What, how, where,</a>	819
763	<a href="#">prehensive survey</a> . <i>Preprint</i> , arXiv:2411.04890.	<a href="#">and how well?</a> <i>Preprint</i> , arXiv:2503.24235.	820
764	Xuehui Wang, Zhenyu Wu, JingJing Xie, Zichen Ding,	Yuqi Zhou, Sunhao Dai, Shuai Wang, Kaiwen Zhou,	821
765	Bowen Yang, Zehao Li, Zhaoyang Liu, Qingyun	Qinglin Jia, and Jun Xu. 2025. <a href="#">Gui-g1: Understand-</a>	822
766	Li, Xuan Dong, Zhe Chen, Weiyun Wang, Xiangyu	<a href="#">ing r1-zero-like training for visual grounding in gui</a>	823
767	Zhao, Jixuan Chen, Haodong Duan, Tianbao Xie,	<a href="#">agents</a> . <i>Preprint</i> , arXiv:2505.15810.	824
768	Chenyu Yang, Shiqian Su, Yue Yu, Yuan Huang,		
769	and 9 others. 2025b. <a href="#">Mmbench-gui: Hierarchical</a>		
770	<a href="#">multi-platform evaluation framework for gui agents</a> .		
771	<i>Preprint</i> , arXiv:2507.19478.		

825	<b>A The Use of Large Language Models (LLMs)</b>	
826		
827	We used LLMs to assist with the phrasing and	
828	grammar of the manuscript. The LLMs were used	
829	strictly as a writing aid and did not contribute to	
830	the scientific ideation, methodology, or results pre-	
831	sented in this paper.	
832	<b>B Implementation Details</b>	
833	For the closed-source models, we perform infer-	
834	ence using the <a href="#">OpenRouter platform</a> , and all use	
835	the default provider. All experiments were con-	
836	ducted between August 15 and September 15 on a	
837	machine equipped with two NVIDIA A100 80GB	
838	GPUs and 1000GB of RAM.	
839	<b>B.1 Benchmark Dataset Discussion</b>	
840	The <b>ScreenSpot-Pro</b> benchmark dataset poses	
841	challenges such as diverse icons, layouts, and	
842	application-specific styles, as well as large input	
843	sizes and heterogeneous content, making it a suit-	
844	able resource for evaluating model robustness.	
845	<b>B.2 Main Experiment Details</b>	
846	In the main experimental results shown in Table 1,	
847	we compare our approach with several baselines:	
848	(1) direct inference using the original models and	
849	(2) DiMo-GUI, one of the strongest existing base-	
850	lines, which incorporates the concept of test-time	
851	scaling. To balance performance and computa-	
852	tional efficiency, we set the threshold for hierar-	
853	chical image search to 600 pixels, meaning that the	
854	search terminates when either the image width or	
855	height falls below this threshold.	
856	<b>B.3 Experimental Model Introduction</b>	
857	Here we briefly introduce the four types of mod-	
858	els that are used in the main experiments, as the	
859	general ‘Scanner’ and the specialist ‘Locator’.	
860	<b>B.3.1 ‘Locator’ Models</b>	
861	• OS-Atlas ( <a href="#">Wu et al., 2024</a> ): An open-source	
862	foundational action model series for GUI	
863	agents, trained on 2.3 million cross-platform	
864	screenshots and 13 million UI elements.	
865	• UGround ( <a href="#">Gou et al., 2025</a> ): A universal visual-	
866	only grounding model family that predicts pixel-	
867	level element coordinates using only visual in-	
868	put, trained on 1.3 million screenshots contain-	
869	ing 10 million GUI elements.	
	<b>B.3.2 ‘Scanner’ Models</b>	870
	• Qwen2.5-VL ( <a href="#">Bai et al., 2025</a> ): A recent multi-	871
	modal vision–language model series, available	872
	in multiple sizes, that offers strong visual un-	873
	derstanding.	874
	• Gemini ( <a href="#">Google, 2025</a> ; <a href="#">Comanici et al., 2025</a> ;	875
	<a href="#">Team et al., 2024, 2025b</a> ): Google’s family of	876
	multimodal models, capable of processing text,	877
	images, audio, and code, and designed for broad	878
	AI applications including chat and search.	879
	<b>C Further Related Work</b>	880
	Here, we further discuss related work concerning	881
	the use of test-time scaling in the field of GUI	882
	grounding:	883
	<b>C.1 Test-Time Scaling</b>	884
	Test-time scaling refers to techniques that improve	885
	model performance at inference without modifying	886
	model parameters, typically by increasing compu-	887
	tation or using additional resources ( <a href="#">Muennighoff</a>	888
	<a href="#">et al., 2025</a> ; <a href="#">Snell et al., 2024</a> ; <a href="#">Zhang et al., 2025c</a> ).	889
	In GUI grounding, test-time scaling has been used	890
	to improve localization through action histories,	891
	external knowledge retrieval, zoom-in searches,	892
	and adaptive focus refinement ( <a href="#">Wu et al., 2025a</a> ;	893
	<a href="#">Nguyen, 2025</a> ; <a href="#">Nakano et al., 2022</a> ). These meth-	894
	ods aim for greater accuracy via extended reason-	895
	ing and iterative attention.	896
	<b>D Additional Baseline Performance</b>	897
	Due to the page limitations in the main paper, we	898
	also report the raw performance of several addi-	899
	tional vision-language models on the test bench-	900
	mark, which we list in Table 3.	901
	From the results presented in the table, it is evi-	902
	dent that most state-of-the-art vision-language mod-	903
	els, including those from families such as GPT-5	904
	and Gemini-2.5-Pro, perform poorly on the GUI	905
	grounding benchmark despite their strong general	906
	vision capabilities. This highlights a critical limita-	907
	tion in their ability to handle fine-grained, domain-	908
	specific grounding tasks. Nevertheless, their robust	909
	visual perception suggests that they can still serve	910
	effectively as visual front-ends to parse and under-	911
	stand GUI images. These findings underscore the	912
	importance of fully leveraging the intrinsic capa-	913
	bilities of such models, rather than relying solely	914
	on scaling up data or fine-tuning larger parameter	915
	models.	916

Base Model	Development			Creative			CAD			Scientific			Office			OS			Average		
	text	icon	avg	text	icon	avg	text	icon	avg	text	icon	avg	text	icon	avg	text	icon	avg	text	icon	avg
Gemma-3-27B	0.0	0.0	0.0	2.0	0.0	1.2	1.0	1.6	1.1	3.5	0.0	2.0	1.1	0.0	0.9	0.0	0.0	0.0	1.3	0.2	0.9
Phi-4-Multimodal	0.0	0.7	0.3	1.5	0.0	0.9	0.5	1.6	0.8	2.1	0.0	1.2	1.7	5.7	2.6	0.0	0.0	0.0	1.0	0.8	0.9
SeeClick	0.6	0.0	0.3	1.0	0.0	0.6	2.5	0.0	1.9	3.5	0.0	2.0	1.1	0.0	0.5	2.8	0.0	1.5	1.8	0.0	1.1
Claude Sonnet 4	1.3	3.4	2.3	1.5	0.7	1.2	2.5	0.0	1.9	1.4	1.8	1.6	0.6	1.9	0.9	0.0	0.0	0.0	1.3	1.5	1.4
Qwen2-VL-7B	2.6	0.0	1.3	1.5	0.0	0.9	0.5	0.0	0.4	6.3	0.0	3.5	3.4	1.9	3.0	0.9	0.0	0.5	2.5	0.2	1.6
GLM-4.5V	0.0	1.4	0.7	1.5	2.1	1.8	5.6	0.0	4.2	2.8	1.0	2.0	1.7	1.9	1.7	0.0	0.0	0.0	2.1	1.2	1.8
GPT-5	2.6	0.7	1.7	4.5	4.2	4.4	7.6	7.8	7.7	4.2	1.8	3.1	4.5	3.8	4.3	0.0	0.0	0.0	4.3	2.6	3.7
Gemini-2.5-Pro	4.5	2.8	3.7	7.6	5.6	6.7	14.2	1.6	11.1	4.9	6.4	5.5	7.3	3.8	6.5	2.8	1.1	2.0	7.5	3.8	6.1
ShowUI-2B	16.9	1.4	9.4	9.1	0.0	5.3	2.5	0.0	1.9	13.2	7.3	10.6	15.3	7.5	13.5	10.3	2.2	6.6	10.8	2.6	7.7
CogAgent-18B	14.9	0.7	8.0	9.6	0.0	5.6	7.1	3.1	6.1	22.2	1.8	13.4	13.0	0.0	6.5	5.6	0.0	3.1	12.0	0.8	7.7
Aria-UI	16.2	0.0	8.4	23.7	2.1	14.7	7.6	1.6	6.1	27.1	6.4	18.1	20.3	1.9	16.1	4.7	0.0	2.6	17.1	2.0	11.3

Table 3: Additional baseline results of various vision-language models on the GUI grounding benchmark. Despite their strong general capabilities, these models perform poorly on this specific task.

## E Baseline Introduction

We compare our GMS framework with several baselines, as shown in Table 1 and Table 3. Below, we introduce each baseline to provide clarification.

- GPT-4o* (OpenAI et al., 2024): OpenAI’s flagship multimodal model that seamlessly understands and generates text, images, and audio. It enables faster, more natural real-time interactions while maintaining strong reasoning and accuracy.
- Gemma3-27B* (Team et al., 2025c): Google’s 27B-parameter version of their Gemma 3 model family. It’s a high-capacity, multimodal model that accepts both text and image inputs, supports an expanded 128K context window, works across 140 languages.
- Phi-4-Multimodal* (Abdin et al., 2024): Microsoft’s 5.6B-parameter model that can process text, vision, and speech (audio) inputs in a unified system. It supports a long 128K token context, uses a ‘mixture of LoRAs’ approach for modality-adapters.
- SeeClick* (Cheng et al., 2024): A visual GUI agent that automates tasks like clicking or typing by observing only interface screenshots, without needing structured representations such as HTML or accessibility trees.
- Claude-Sonnet-4* (Anthropic, 2025): A mid-tier model in Anthropic’s Claude 4 family, designed to balance strong reasoning and coding ability with efficiency and accessibility.
- Qwen2-VL-7B* (Wang et al., 2024b; Bai et al., 2023): A 7B-parameter vision-language model from Alibaba’s Qwen2-VL family.
- GLM-4.5V* (Team et al., 2025a): ZhipuAI’s flagship vision-language model built on GLM-4.5-Air, activating 12B of its 106B parameters per pass to balance efficiency with strong multimodal reasoning.
- Gemini-2.0-Flash* (Google, 2025): Google’s high-performance, multimodal model in the Gemini 2.0 family designed for the ‘agentic era’.
- Gemini-2.5-Flash-Lite* (Comanici et al., 2025): Google’s cost- and latency-optimized variant in the Gemini 2.5 model series, designed for high-volume, real-world use.
- GPT-5* (OpenAI, 2025): OpenAI’s next-generation multimodal model that advances beyond GPT-4o with stronger reasoning, longer context handling, and more efficient real-time interaction across text, vision, and audio.
- Gemini-2.5-Pro* (Comanici et al., 2025): Google’s top-tier reasoning model in the Gemini 2.5 family, designed to tackle complex problems across modalities, including text, audio, images, video, and even whole code repositories.
- ShowUI-2B* (Lin et al., 2024): A lightweight vision-language-action model from ShowLab, built for GUI agents to understand and interact with graphical user interfaces via screenshots.

979  
980  
981  
982  
983  
  
984  
985  
986  
987  
988  
  
989  
990  
991  
992  
993  
  
994  
995  
996  
997  
998  
  
999  
1000  
  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
  
1012  
  
1013  
1014  
1015  
1016  
1017  
  
1018  
  
1019  
1020  
1021  
1022  
1023  
1024

- *CogAgent-18B* (Hong et al., 2024): An open-source vision-language model (VLM) developed by THU DM and Zhipu AI, specifically optimized for understanding and interacting with graphical user interfaces (GUIs).
- *Aria-UI* (Yang et al., 2025): A multimodal model for GUI grounding that maps language instructions to specific interface elements using only vision (screenshots), foregoing HTML or accessibility trees (AXTrees) as auxiliary input.
- *Claude (Computer Use)* (Hu et al., 2024): A GUI-agent extension of Claude 3.5 Sonnet that enables the model to observe screenshots of a user’s computer and issue desktop actions (mouse, keyboard, clicks) to automate tasks.
- *UI-TARS-7B* (Qin et al., 2025): A 7B-parameter vision-language model from ByteDance designed for native GUI automation, capable of controlling both web and desktop applications via only screenshot input.
- *UI-TARS-72B* (Qin et al., 2025): 72B-parameter version of UI-TARS.

For the baselines not presented in the previous paper, we conduct the evaluations ourselves, with the inference prompts provided in Appendix G. For the baselines included in the previous paper, we directly use the results reported by Wu et al. (2025a), which also correspond to the ScreenSpot-Pro leaderboard data. Regarding the general-purpose vision-language models, we select recent and strong models to demonstrate their raw performance on the direct GUI grounding task, which turns out to be rather poor.

## F Discussion

In this section, we elaborate on the proposed framework and present further analysis of the associated experiments, demonstrating its superior performance and the novel insights it yields relative to existing methods.

### F.1 Cognition and Architectural Insights

Here, we discuss several key aspects in which our GMS framework departs from prior works, highlighting the unique insights underlying our design:

- (1) Cognitive-inspired task decomposition based on the dual-stream model, separating semantic attention from motor-level localization.

- (2) Hierarchical attention and cross-modal verification that iteratively refine the search space, replacing brittle single-pass grounding.
  - (3) Asymmetric multi-agent collaboration, with a generalist scanner for abstraction and a specialist locator for spatial precision.
  - (4) Late-stage fusion through multi-resolution decision making, aligning global predictions with fine-grained local cues.
- These insights allow GMS to achieve a stronger balance between global semantic reasoning and local spatial precision than prior approaches.

### F.2 On the Possibility of Comparing with Additional Baselines

In this paper, we primarily compare our framework with DiMo-GUI, one of the strongest existing baselines on the GUI grounding benchmark. Although numerous related works report results on this benchmark, their performance under comparable grounding model settings consistently falls short of DiMo-GUI. Therefore, we focus our comparison on DiMo-GUI to balance both reproducibility costs and page constraints. Given that our framework outperforms DiMo-GUI, it is reasonable to infer that it also surpasses other baselines that perform worse than DiMo-GUI.

## G Inference Prompts

We present the manually designed inference prompts employed in the experiments shown in Figures 6 through 18.

The first thing we want to note is that, for the specific fine-tuned grounding models, we use the same inference prompt as the researchers in the previous work, without making any modifications. The inference-related code is written based on the HuggingFace repository example code, including some resizing and transformations for certain models such as OS-Atlas-4B.

The second point to note is that each main prompt designed for our GMS framework has two versions, depending on the names of certain subsets. Subsets such as "ppt\_windows" or "word\_macos" clearly indicate the application name (here "powerpoint" and "word"). However, there are three special subsets, namely "common\_linux", "common\_windows", and "common\_macos", which do not contain specific application names. For this reason, we provide two versions of each prompt: the first is used for most subsets, while the second is

1025  
1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
  
1037  
1038  
  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
  
1051  
  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073

1074 used for the three special subsets mentioned here.

**Hierarchical Attention Allocation Prompt (Initial Level & Normal Version)**

I have provided you a screenshot of my desktop containing the interface of the {application\_name} application running on the {system\_name} system. Where should I click if I want to DIRECTLY perform the following operation in the {application\_name}: \*\*{instruction}\*\*? Provide the possibilities for each region (Region 1 to Region 9, ordered from left to right, top to bottom) with a score between 0 and 100. Your output MUST follow this format: "Region X: SCORE (explanation)".

Figure 6: The hierarchical attention allocation prompt for the initial level (search depth = 0) and the normal subsets.

**Hierarchical Attention Allocation Prompt (Initial Level & Special Version)**

I have provided you a screenshot of my desktop using {system\_name} system. Where should I click if I want to DIRECTLY perform the following operation in the {application\_name}: \*\*{instruction}\*\*? Provide the possibilities for each region (Region 1 to Region 9, ordered from left to right, top to bottom) with a score between 0 and 100. Your output MUST follow this format: "Region X: SCORE (explanation)".

Figure 7: The hierarchical attention allocation prompt for the initial level (search depth = 0) and the special subsets.

**Hierarchical Attention Allocation Prompt (Non-Initial Level & Normal Version)**

Where should I click if I want to DIRECTLY perform the following operation in the {application\_name}: \*\*{instruction}\*\*? Provide the possibilities for each region (Region 1 to Region 9, ordered from left to right, top to bottom) with a score between 0 and 100. Your output MUST follow this format: "Region X: SCORE (explanation)".

Figure 8: The hierarchical attention allocation prompt for the non-initial level (search depth  $\geq 1$ ) and the normal subsets.

**Hierarchical Attention Allocation Prompt (Non-Initial Level & Special Version)**

Where should I click if I want to DIRECTLY perform the following operation: \*\*{instruction}\*\*? Provide the possibilities for each region (Region 1 to Region 9, ordered from left to right, top to bottom) with a score between 0 and 100. Your output MUST follow this format: "Region X: SCORE (explanation)".

Figure 9: The hierarchical attention allocation prompt for the non-initial level (search depth  $\geq 1$ ) and the special subsets.

**Scanner Region Verification Prompt (Normal Version)**

You need to check if the image region from my desktop screenshot contains the button or icon for me to DIRECTLY perform the following operation in the {application}: \*\*{instruction}\*\*. You are required to output your reasoning process first, and then provide your final answer in the format: <answer>yes/no</answer>.

Figure 10: The region verification prompt that instructed the 'Scanner' agent to filter the region of interest (for normal subsets).

#### Scanner Region Verification Prompt (Special Version)

You need to check if the image region from my desktop screenshot contains the button or icon for me to DIRECTLY perform the following operation: **{instruction}**. You are required to output your reasoning process first, and then provide your final answer in the format: `<answer>yes/no</answer>`.

Figure 11: The region verification prompt that instructed the ‘Scanner’ agent to filter the region of interest (for special subsets).

#### Cross-Modal Verification Prompt (Normal Version)

I want to DIRECTLY perform the following operation in the {application}: **{instruction}**. First, I’m showing you the original screenshot image, followed by a cropped region from it. You need to determine whether this cropped region is likely to contain the target button, area, or icon for the operation. Answer with `<relevance>yes/no</relevance>` and provide your reasoning within `<reasoning>...</reasoning>`.  
This is the original screenshot image: `<Image1>`  
And this is the cropped region from the original screenshot image: `<Image2>`

Figure 12: The cross-modal verification prompt that instructed the ‘Scanner’ agent to verify the cropped region (for normal subsets).

#### Cross-Modal Verification Prompt (Special Version)

I want to DIRECTLY perform the following operation: **{instruction}**. First, I’m showing you the original screenshot image, followed by a cropped region from it. You need to determine whether this cropped region is likely to contain the target button, area, or icon for the operation. Answer with `<relevance>yes/no</relevance>` and provide your reasoning within `<reasoning>...</reasoning>`.  
This is the original screenshot image: `<Image1>`  
And this is the cropped region from the original screenshot image: `<Image2>`

Figure 13: The cross-modal verification prompt that instructed the ‘Scanner’ agent to verify the cropped region (for special subsets).

#### Scanner Adaptive Resolution Enhancement Prompt (Normal Version)

I want to DIRECTLY perform this operation in the {application} on my desktop: **{instruction}**.  
I have extracted a candidate region and divided it into 5x5 smaller regions (numbered 1 to 25 from left to right, top to bottom). Please identify which of the 25 regions is the most relevant (return only one region you are most confident about). Then, within that region, tell me which of the 9 inner zones the target click point is closest to. (Choose from: top left, top center, top right, center left, center, center right, bottom left, bottom center, bottom right)  
First, provide your reasoning process, and then return your final answer in the following format:  
`<index>xxx</index>`  
`<location>xxx</location>`

Figure 14: The adaptive resolution prompt that instructed the ‘Scanner’ agent to identify a possible fine-grained region (for normal subsets).

#### Scanner Adaptive Resolution Enhancement Prompt (Special Version)

I want to DIRECTLY perform this operation on my desktop: **\*\*{instruction}\*\***.  
I have extracted a candidate region and divided it into 5x5 smaller regions (numbered 1 to 25 from left to right, top to bottom). Please identify which of the 25 regions is the most relevant (return only one region you are most confident about). Then, within that region, tell me which of the 9 inner zones the target click point is closest to. (Choose from: top left, top center, top right, center left, center, center right, bottom left, bottom center, bottom right)  
First, provide your reasoning process, and then return your final answer in the following format:  
<index>xxx</index>  
<location>xxx</location>

Figure 15: The adaptive resolution prompt that instructed the ‘Scanner’ agent to identify a possible fine-grained region (for special subsets).

#### OS-Atlas-4B Grounding Prompt

In the screenshot of this web page, please give me the coordinates of the element I want to click on according to my instructions(with point).\n“{ }”

Figure 16: The instruction for the OS-Atlas-4B model to output grounding coordinates.

#### UGround-7B Grounding Prompt

In the screenshot, where are the pixel coordinates (x, y) of the element corresponding to “{ }”?

Figure 17: The instruction for the UGround-7B model to output grounding coordinates.

#### UGround-V1-7B Grounding Prompt

Your task is to help the user identify the precise coordinates (x, y) of a specific area/element/object on the screen based on a description.

- Your response should aim to point to the center or a representative point within the described area/element/object as accurately as possible.
- If the description is unclear or ambiguous, infer the most relevant area or element based on its likely context or purpose.
- Your answer should be a single string (x, y) corresponding to the point of the interest.

Description: {instruction}  
Answer:

Figure 18: The instruction for the UGround-V1-7B model to output grounding coordinates.

#### Baseline Models Grounding Prompt

I want to DIRECTLY perform this operation in the {application} on my desktop: **\*\*{instruction}\*\***. You should provide the target CLICK pixel coordinate (x, y) in the ORIGINAL image. You must output only integer coordinate values. For example: '123, 456' or '(123, 456)'.

Figure 19: The instruction for the baseline models to output grounding coordinates.