# InfiGUIAgent: A Multimodal Generalist GUI Agent with Native Reasoning and Reflection

Yuhang Liu<sup>1</sup> Pengxiang Li<sup>2</sup> Zishu Wei<sup>1</sup> Congkai Xie<sup>3</sup> Xueyu Hu<sup>1</sup> Xinchen Xu<sup>1</sup> Shengyu Zhang<sup>1</sup> Xiaotian Han<sup>4</sup> Hongxia Yang<sup>5</sup> Fei Wu<sup>1</sup>

### Abstract

Graphical User Interface (GUI) Agents, powered by multimodal large language models (MLLMs), have shown great potential for task automation on computing devices such as computers and mobile phones. However, existing agents face challenges in multi-step reasoning and reliance on textual annotations, limiting their effectiveness. We introduce InfiGUIAgent, an MLLM-based GUI Agent trained with a two-stage supervised finetuning pipeline. Stage 1 enhances fundamental skills such as GUI understanding and grounding, while Stage 2 integrates hierarchical reasoning and expectation-reflection reasoning skills using synthesized data to enable native reasoning abilities of the agents. InfiGUIAgent achieves competitive performance on several GUI benchmarks, highlighting the impact of native reasoning skills in enhancing GUI interaction for automation tasks.

## 1. Introduction

Graphical User Interface (GUI) Agents have emerged as powerful tools for automating tasks on computing devices, including mobile phones and computers. These agents can understand and interact with GUIs to execute complex operations, significantly enhancing user productivity and expanding the scope of automated task completion (Hu et al., 2024b; Hong et al., 2024; Zhang & Zhang, 2023; Qi et al., 2024; Xie et al., 2024; Vu et al., 2024; Yu et al., 2023). Recent developments in multimodal large language models (MLLMs) (Bai et al., 2023b; Li et al., 2024c; Team et al., 2024; Dai et al., 2022) have significantly advanced the potential of GUI Agents. MLLMs possess



*Figure 1. InfiGUIAgent*, trained via a two-stage SFT pipeline, addresses key GUI agent challenges. (a) Baseline grounding is often unreliable, with auxiliary text offering partial gains at the cost of overhead, while pure vision struggles. (b) Stage 1 enhances fundamental GUI understanding and grounding, but agents remain prone to repetitive errors in multi-step trajectory tasks. (c) Stage 2 integrates native hierarchical and expectation-reflection reasoning, empowering *InfiGUIAgent* for successful complex task automation.

powerful visual understanding capabilities and can reason based on visual information, making them a promising foundation for building sophisticated GUI Agents. These models can interpret complex interface elements and adapt to a wide range of tasks, leading to more efficient and robust automation (Hong et al., 2024; Jiang et al., 2023; You et al., 2025; Nong et al., 2024; Vu et al., 2024).

Despite their promise, MLLM-based GUI Agents encounter hurdles in achieving robust, autonomous task completion. A primary challenge lies in their interaction with the visual interface. Many agents struggle to connect natural language instructions to the correct GUI elements (such as icons or buttons) within complex layouts. Determining the precise location for an interaction, like a tap, also remains unreliable. This often compels a reliance on additional in-

<sup>&</sup>lt;sup>1</sup>Zhejiang University <sup>2</sup>Dalian University of Technology <sup>3</sup>Reallm Labs <sup>4</sup>Independent <sup>5</sup>The Hong Kong Polytechnic University. Correspondence to: Shengyu Zhang <sy\_zhang@zju.edu.cn>.

*Workshop on Computer-use Agents @ ICML 2025, Vancouver, Canada. Copyright 2025 by the author(s).* 

formation, such as accessibility trees or Set-of-Marks (Yang et al., 2023b), to interpret the GUI. However, GUIs are inherently visual; textual augmentation can lead to information loss, increased computational overhead, and platform inconsistencies, hindering practical deployment.

Beyond these interaction challenges, agents face limitations in their reasoning and planning capabilities (Zhang & Zhang, 2023; Qi et al., 2024; Yu et al., 2024). Many do not effectively leverage insights from previous steps or reflect on past actions, frequently leading to repetitive errors. Separately, agents often make decisions that yield low long-term value or are incorrect, indicating difficulties in decomposing high-level goals into effective sub-steps. These issues, particularly acute in MLLMs with fewer parameters, can cause agents to become trapped in unproductive behavioral loops.

Such behaviors suggest underlying weaknesses: a deficit in robust visual-semantic grounding and spatial awareness at the perceptual level, and at the cognitive level, a lack of intrinsic mechanisms for proactive historical analysis (hindering self-correction) and strategic, multi-scale planning.

Overcoming these limitations requires more than applying existing MLLMs or relying on elaborate prompting. We introduce InfiGUIAgent, an MLLM-based GUI Agent developed through a two-stage supervised fine-tuning (SFT) pipeline (Figure 1). This pipeline is architected to: first, instill foundational GUI understanding and grounding capabilities (Stage 1) by utilizing diverse datasets covering layout comprehension, grounding tasks, and GUI QA, thereby directly improving the agent's ability to interpret and interact with visual interfaces. Second, cultivate advanced reasoning skills (Stage 2)-hierarchical reasoning and expectation-reflection reasoning-using carefully synthesized data from trajectories, enabling these to be performed 'natively' by the agent. This overall approach aims to equip the agent with internalized mechanisms for planning and adaptive learning from its operational history. Infi-GUIAgent achieves competitive performance on several GUI benchmarks, highlighting the impact of this integrated approach. Our main contributions are threefold:

- We propose a two-stage supervised fine-tuning pipeline to comprehensively improve both the fundamental abilities and advanced reasoning abilities of GUI Agents.
- We synthesize SFT data with two advanced reasoning skills: hierarchical reasoning and expectation-reflection reasoning, enabling the agents to natively perform complex reasoning. This data will be open-sourced to promote community development.
- We build *InfiGUIAgent* by supervised fine-tuning a model using our SFT data and conduct experiments on several GUI benchmarks, demonstrating that our model achieves competitive performance.

#### 2. Related Works

#### 2.1. Multimodal LLMs

Large Language Models (LLMs) (Floridi & Chiriatti, 2020; Touvron et al., 2023; Bai et al., 2023a; Xiao et al., 2021) have significantly enhanced the capabilities of AI systems in tackling a wide range of tasks (Hu et al., 2024c; Li et al., 2024d), thanks to their exceptional ability to process complex semantic and contextual information. The remarkable power of LLMs has also inspired exploration into their potential for processing multimodal data, such as images. Typically, the architecture of Multimodal Large Language Models (MLLMs) consists of three main components: a pre-trained large language model, a trained modality encoder, and a modality interface that connects the LLM with the encoded modality features. Various vision encoders, such as ViT (Dosovitskiy et al., 2021), CLIP (Radford et al., 2021), and ConvNeXt (Liu et al., 2022), extract visual features, which are integrated using techniques like adapter networks (Liu et al., 2023), cross-attention layers (Alayrac et al., 2022), and visual expert modules (Wang et al., 2023). These methods have facilitated the development of highperforming MLLMs, such as Qwen-VL (Bai et al., 2023b), GPT-4 Vision (OpenAI, 2023), BLIP-2 (Li et al., 2023) and InfiMM (Liu et al., 2024), thus opening new avenues for LLMs in processing GUI tasks.

#### 2.2. MLLM-based GUI Agents

Agents are AI systems that perceive their environments, make decisions, and take actions to complete specific tasks. LLMs reaching human-level intelligence have greatly enhanced the ability to build agents. For GUI tasks, LLMs that read HTML code to perceive GUIs are developed (Wen et al., 2023). However, various works have shown that learning to interact with the visual form of the GUIs can show superior performance (Hu et al., 2024b). Therefore, MLLM-based GUI Agents are developed. ILuvUI (Jiang et al., 2023) fine-tuned LLaVA to enhance general GUI understanding, while AppAgent (Zhang et al., 2023) explored app usage through autonomous interactions. CogAgent (Hong et al., 2024) integrated high-resolution vision encoders, and Ferret-UI-anyres (You et al., 2025) employed an any-resolution approach. Building upon these works, our study focuses on developing a more lightweight agent with a simplified architecture for GUI tasks, aiming to improve ease of deployment.

## 3. Method

In this section, we introduce our two-stage supervised finetuning strategy for building *InfiGUIAgent*, as shown in Figure 2. In stage 1, we focus on improving fundamental abilities such as understanding and grounding, particularly



#### **Stage 1: Fundamental Abilities**

#### Stage 2: Native Advanced Reasoning

Figure 2. InfiGUIAgent is trained in two stages. Stage 1 cultivates fundamental abilities using diverse datasets covering GUI understanding (element recognition and layout comprehension), question answering, instruction grounding, general knowledge, and tool usage. Stage 2 introduces native advanced reasoning, employed during both training and inference. This stage follows a cyclical process at each step, consisting of Reflection, Hierarchical Reasoning (strategic and tactical layers), Action, and Expectation. Each step receives the overall task, the history of previous screenshots and reasoning, and the current environment as input. Reflection assesses the previous action's outcome against its expectation, while Expectation predicts the outcome of the current action for subsequent reflection.

considering the complexity of GUIs. In stage 2, we move on to improve the native reasoning abilities of agents for handling complicated GUI tasks.

#### 3.1. Stage 1: Training for Fundamental Abilities

Considering the complexity of GUIs, which involve diverse data formats such as HTML code, high-resolution interfaces cluttered with small icons and text, general MLLMs lack fundamental abilities in both understanding GUI and grounding the actions. To address this, we first collected a range of existing visual-language and GUI datasets for supervised fine-tuning in stage 1. We gathered data covering several GUI tasks from multiple sources to ensure a comprehensive capabilities improvement (see Table 1). The datasets can be categorized into five parts:

- **GUI Understanding.** Datasets focusing on GUI element recognition, layout comprehension, and semantic interpretation, including Screen2Words (Wang et al., 2021) and Screen Annotation (Baechler et al., 2024).
- Grounding. Datasets capture various user interaction sequences and operation patterns, including GUIEnv (Chen et al., 2024), RICO Semantic Annotation (Sunkara et al., 2022), SeeClick-Web (Cheng et al., 2024), RICO SCA (Li et al., 2020a), Widget Caption (Li et al., 2020b), UIBert Reference Expression (Bai et al., 2021) and OmniAct-Single Click (Kapoor et al., 2024).
- Question Answering. Datasets contain GUI-specific QA

tasks, including GUIChat (Chen et al., 2024), ScreenQA (Hsiao et al., 2022) and Complex QA (Yin et al., 2023).

- General Knowledge. Multimodal datasets maintain model's general capabilities, including LLaVA-OneVision (Li et al., 2024b) and PixMo (MDeitke et al., 2024).
- **Tool Usage.** Datasets cover general tool using, including Glaive-function-calling (Glaive AI, 2024).

Due to the diversity of our data sources, we implemented comprehensive format standardization across all datasets. Additionally, we adopted the Reference-Augmented Annotation format (see Section 3.1.2) to enhance the model's ability to ground visual elements with textual descriptions, enabling precise spatial referencing while maintaining natural language flow.

#### 3.1.1. DATA PREPROCESSING AND STANDARDIZATION

Given the diversity of our data sources, we implemented comprehensive preprocessing steps to standardize the data format across all datasets. We normalized the coordinate system by following (Wang et al., 2024), mapping all spatial coordinates to a relative scale of [0, 1000]. This standardization facilitates consistent representation of both point and box annotations in JSON format, with points expressed as  $\{"x": x, "y": y\}$  and bounding boxes as  $\{"x1": x_1, "y1": y_1, "x2": x_2, "y2": y_2\}$ . In this coordinate system, the origin  $\{"x": 0, "y": 0\}$  is located at the screen's top-left corner, with the x-axis extending rightward

Dataset	Platform	Category	# of Samples
GUI-related Datasets			
GUIEnv (Chen et al., 2024)	Webpage	Grounding	150,000
RICO Semantic Annotation (Sunkara et al., 2022)	Mobile	Grounding	150,000
SeeClick-Web (Cheng et al., 2024)	Webpage	Grounding	100,000
RICO SCA (Li et al., 2020a)	Mobile	Grounding	100,000
Widget Caption (Li et al., 2020b)	Mobile	Grounding	70,000
GUIChat (Chen et al., 2024)	Webpage	QA	40,000
ScreenQA (Hsiao et al., 2022)	Mobile	QA	17,000
UIBert Reference Expression (Bai et al., 2021)	Mobile & Mobile	Grounding	16,000
Screen2Words (Wang et al., 2021)	Mobile	Understanding	12,000
Complex QA (Yin et al., 2023)	Mobile	QA	11,000
Screen Annotation (Baechler et al., 2024)	Mobile	Understanding	5,400
OmniAct-Single Click (Kapoor et al., 2024)	Webpage & Desktop	Grounding	4,800
Non-GUI Datasets			
LLaVA-OneVision (Li et al., 2024b)	-	General	250,000
PixMo (MDeitke et al., 2024)	-	General	68,800
Glaive-function-calling (Glaive AI, 2024)	-	Tool Usage	5,000

Table 1. Training datasets used in stage 1 of supervised fine-tuning.

and the y-axis downward. The bottom-right corner corresponds to coordinates {"x" : 1000, "y" : 1000}. To enhance data quality, we implemented two additional preprocessing steps:

**Instruction Enhancement.** For datasets with ambiguous instructions (e.g., those varying in phrasing or lacking explicit intent across different data sources), we developed standardized instruction templates to establish clear correspondence between commands and their expected outcomes.

**Response Refinement.** For entries with complex or inconsistent response formats, we utilized Qwen2-VL-72B (Bai et al., 2023b) to reformulate responses while preserving their semantic content. Each reformulation underwent validation to ensure accuracy and consistency.

#### **3.1.2. Reference-Augmented Annotation**

To better leverage the spatial information available in our collected datasets and enhance the model's visual-language understanding of GUIs, we implemented a referenceaugmented annotation format. This format enables bidirectional referencing between GUI elements and textual responses. Specifically, we adopted the following structured notation:

The format consists of several key components: the refer-

ence type (either "box" for rectangular regions or "point" for specific locations), coordinate specifications (x1, y1, x2, y2 for boxes or x, y for points), optional annotative notes, and the corresponding textual content. To generate training data in this format, we prompted Qwen2-VL-72B (Bai et al., 2023b) to seamlessly integrate GUI spatial information with original responses, maintaining natural language flow while preserving precise spatial references.

#### 3.2. Stage 2: Training for Native Reasoning

Building upon the foundational capabilities such as understanding and grounding, GUI Agents must also master advanced reasoning skills to effectively handle complex tasks. We identify two crucial reasoning skills : (1) Hierarchical reasoning, which enables planning and task decomposition, helping agents structure complex tasks into manageable subtasks and execute them efficiently (Huang & Chang, 2023; Zhang et al., 2024b; Huang et al., 2024), and (2) Expectation-reflection reasoning, which fosters adaptive self-correction and reflection (Shinn et al., 2023; Yao et al., 2023; Hu et al., 2024a), enabling agents to learn from past actions and improve decision-making consistency. These reasoning skills are integrated into the training datasets of agents, so that they can reason with these skills natively without any extra prompting. To achieve this, we generate SFT data incorporating these reasoning skills based on existing trajectory data (see Table 2) and continue fine-tuning the model from stage 1.

#### 3.2.1. HIERARCHICAL REASONING

Effective execution of GUI tasks demands both overarching strategic planning and meticulous tactical execution. To achieve this, we synthesize trajectory data with a hierarchi-

Dataset	Platform	# of Samples
GUIAct (Chen et al., 2024)	Webpage & Mobile	10,000
AMEX (Chai et al., 2024)	Mobile	3,000
Android in the Zoo (Zhang et al., 2024a)	Mobile	2,000
Composition: Stage 1-aligned	-	30,000

Category	Operations
Single-point operations	tap, click, hover, select
Two-point operations	swipe, select_text
Directional operations	scroll
Text input	input,point_input
Parameterless operations	remember, enter, home, back
State settings	set_task_status

*Table 3.* Categorization of actions in the action space.

cal reasoning with two distinct layers:

- **Strategic Layer.** Strategic layer is responsible for highlevel task decomposition and sub-goal planning. This layer analyzes the overall task objective and determines the sequence of subtasks needed for completion.
- **Tactical Layer.** Tactical layer handles the selection and grounding of concrete actions. Based on the strategic layer's planning, agent select appropriate GUI operations and adjusts their parameters to match the target.

#### 3.2.2. EXPECTATION-REFLECTION REASONING

To enhance action consistency and foster autonomous selfcorrection, we incorporate Expectation-reflection reasoning into the training datasets. This iterative process enhances the agent's ability to adapt and learn from its actions through a structured reflection cycle:

- **Reasoning.** After reflection (except the first step), the agents conduct hierarchical reasoning.
- Action. After the reasoning, the agent takes the action.
- **Expectation.** Following each action, the agent generates expected outcomes which are used to be verified at the next step.
- **Reflection.** The agent evaluates whether its actions achieved the expected results and generating a textual summary of the reflection.

#### 3.2.3. Agent-Environment Interface

We formulate the GUI interaction as a process where an agent interacts with a mobile environment. Let  $s_t \in S$  denote the environment state at step t, where S represents the state space. The agent can observe the state through

a screenshot observation  $o_t$  and performs actions  $a_t \in \mathcal{A}$ , where  $\mathcal{A}$  is the action space. The environment transitions from  $s_t$  to  $s_{t+1}$  following  $s_{t+1} \sim P(\cdot|s_t, a_t)$ , where Prepresents the transition probability function.

The agent receives a task goal g and maintains access to a history window of size n. At each step t, the agent's input consists of:

- Goal g
- Current observation  $o_t$
- Historical context  $H_t = \{(o_i, r_i, a_i)\}_{i=t-n}^{t-1}$ , where  $r_i$  represents the reasoning process

Based on these inputs, the agent generates a reasoning process  $r_t$  and predicts an action  $a_t$ . The interaction follows a standard protocol using function calls and responses:

```
Assistant Message:
<tool_call>
{
    "name": "action_name",
    "arguments": {"action_parameters"}
}
</tool_call>
Tool Message:
<tool_response>
{
    "name": "gui_operation",
    "content": {
       "status": "success | failure",
       "current_ui": <image>,
       "current_task": <task_description>
   }
}
</tool_response>
```

### 3.2.4. MODULAR ACTION SPACE

Given the diverse action spaces across collected datasets, we categorized and standardized the actions by unifying their names and parameters, merging similar operations where appropriate. The resulting action space  $\mathcal{A}$  consists of independent, composable operations that can be flexibly combined based on task requirements, as shown in Table 3.

Table 2. UI action reasoning	datasets	used in th	e training	process
------------------------------	----------	------------	------------	---------

This modular design allows for dynamic action space configuration while maintaining a consistent interface across different platforms and scenarios. The modularity of the action space is implemented by modifying the system prompt. The system prompt we use is shown in Appendix B.2.

#### 3.2.5. REASONING PROCESS CONSTRUCTION

To construct high-quality reasoning data to stimulate the model's native reasoning capabilities, we leverage more capable MLLMs (e.g. Qwen2-VL-72B) to generate structured reasoning processes based on existing interaction trajectories. The construction process involves several key components:

- Screenshot Description. For each observation  $o_t$  in the trajectory, we generate a detailed description  $d_t$ . This step addresses the limitation that some MLLM models do not support interleaved image-text input formats well. To establish clear correspondence between observations (screenshots) and steps, we generate detailed descriptions to replace the screenshots, which helps facilitate the subsequent reasoning process construction.
- **Reflection.** Given the previous expectation  $e_{t-1}$  and current observation  $o_t$ , we generate a reflection  $f_t$  that evaluates the outcome of the previous action.
- Strategic Layer. The strategic reasoning consists of two parts: First, a summary is generated based on the n-step history  $H_t = \{(o_i, r_i, a_i)\}_{i=t-n}^{t-1}$  and current observation  $o_t$ . Then, the planning component is generated with access to the actual action  $a_t$  to ensure alignment with the trajectory.
- Tactical Layer. This layer's reasoning is constructed using the generated reflection  $f_t$  and strategic layer output. The actual action  $a_t$  from the trajectory is incorporated to ensure the tactical reasoning leads to appropriate action selection.
- Expectation. For each state-action pair  $(s_t, a_t)$ , we generate an expectation  $e_t$  based on current observation  $o_t$ , reasoning process  $r_t$ , and action  $a_t$ . Notably, we deliberately avoid using the next state  $s_{t+1}$  in this generation process. Although using  $s_{t+1}$  could improve the agent's accuracy in modeling state transitions, while using  $s_{t+1}$  could lead to perfect expectations, such an approach might impair the agent's ability to handle expectation mismatches during deployment.

While we avoid using  $s_{t+1}$  in expectation generation to maintain robustness, we also explore the possibility of improving state transition modeling through a parallel nextstate prediction task. Using the trajectory data, we construct additional training examples where the agent learns to predict the next state description  $d_{t+1}$  given the current observation  $o_t$  and action  $a_t$ . This auxiliary task helps the agent learn state transition dynamics, while keeping the expectation generation process independent of future states.

## 4. Experiments

#### 4.1. Experimental Setting

#### 4.1.1. IMPLEMENTATION DETAILS

In stage 1, we sample 1M samples in total as illustrated in Table 1. In stage 2, we synthesized 45K samples based on trajectories from datasets shown in Table 2. We employ a full fine-tuning approach to continually supervised fine-tune Qwen2-VL-2B (Bai et al., 2023c). The hyperparameters used for training were a learning rate of  $5 \times 10^{-6}$ , a batch size of 256, and a maximum sequence length of 32k to accommodate long trajectories. The training was conducted for 1 epoch with a warmup ratio of 0.05. To optimize the learning process across the two stages, the vision module was unfrozen during the first stage to facilitate the learning of fundamental UI knowledge, and subsequently frozen during the second stage to allow the model to focus more on reasoning capabilities. We leverage ZeRO (stage 0) (Rajbhandari et al., 2020) technology and FlashAttention-2 (Dao, 2023) to accelerate training and reduce memory consumption, enabling full parameter fine-tuning of the model across 8 A800 80GB GPUs.

#### 4.1.2. EVALUATION BENCHMARKS

**ScreenSpot.** ScreenSpot (Cheng et al., 2024) is an evaluation benchmark for GUI grounding, consisting of over 1,200 instructions from iOS, Android, macOS, Windows, and Web environments, with annotated element types.

AndroidWorld. AndroidWorld (Rawles et al., 2024) is a fully functional Android environment that provides reward signals for 116 programmatic tasks across 20 real-world Android apps. We find that Android World uses Set-of-Marks (SoM) (Yang et al., 2023a) to enhance the agent's grounding ability. However, when humans operate smartphones, their brains do not label elements on the screen. Over-reliance on SoM can lead to insufficient focus on pixel-level grounding ability. Therefore, in our experiments, agents respond to the raw image rather than the annotated image.

## 4.2. Main Results

**ScreenSpot.** Table 4 provides the results of different models across three platforms (Mobile, Desktop and Web) and two element types (Text and Icon) on ScreenSpot (Cheng et al., 2024). InfiGUIAgent-2B achieves highest accuracy of 76.3%, surpassing several strong baselines such as ShowUI (Lin et al., 2024) (75.1%) and UGround-7B (Gou et al., 2024) (73.3%), which is even with larger parameters size.

InfiGUIAgent: A Multimodal Generalist GUI Agent with Native Reasoning and Reflection

Model	Accuracy (%)           Mobile         Desktop						
			Desktop		Web		
	Text	Icon	Text	Icon	Text	Icon	
Proprietary Models							
GPT-40 <sup>1</sup> (OpenAI, 2024)	30.5	23.2	20.6	19.4	11.1	7.8	18.8
Gemini-1.5-pro <sup>2</sup> (Team et al., 2024)	76.2	54.1	65.5	39.3	52.2	32.0	53.2
Open-source Models							
Qwen2-VL-2B (Wang et al., 2024)	24.2	10.0	1.4	9.3	8.7	2.4	9.3
Qwen2-VL-7B (Wang et al., 2024)	61.3	39.3	52.0	45.0	33.0	21.8	42.9
Qwen2-VL-72B (Wang et al., 2024)	73.6	65.5	54.1	57.1	53.0	44.2	57.9
CogAgent (Hong et al., 2024)	67.0	24.0	74.2	20.0	70.4	28.6	47.4
SeeClick (Cheng et al., 2024)	78.0	52.0	72.2	30.0	55.7	32.5	53.4
UGround-7B (Gou et al., 2024)	82.8	60.3	82.5	<u>63.6</u>	80.4	70.4	73.3
ShowUI-2B (Lin et al., 2024)	92.3	75.5	76.3	61.1	<u>81.7</u>	<u>63.6</u>	<u>75.1</u>
Ours							
InfiGUIAgent-2B	92.3	<u>70.3</u>	85.6	65.7	83.0	<u>63.6</u>	76.8

*Table 4.* Performances on various platforms (Mobile, Desktop, Web) on Screenshot. All experiments were conducted using raw screenshot information. Results marked in **bold** represent the best performance, and those <u>underlined</u> indicate the second-best performance.

InfiGUIAgent-2B.

Model	Success Rate			
	Easy	Middle	Hard	Overall
Open-source Models				
Qwen2-VL-2B (Wang et al., 2024)	0.00	0.00	0.00	0.00
Qwen2-VL-7B (Wang et al., 2024)	0.00	0.00	0.00	0.00
Qwen2-VL-72B (Wang et al., 2024)	0.08	0.00	0.00	0.04
LLaVa-OV-7B (Li et al., 2024a)	0.00	0.00	0.00	0.00
ShowUI-2B (Lin et al., 2024)	0.18	0.00	0.00	0.09
Ours				
InfiGUIAgent-2B	0.25	0.00	0.00	0.13

Configuration	ScreenSpot (% Accuracy)	AndroidWorld (Success Rate)
InfiGUIAgent-2B	76.8	0.13
Ablations		
w/o Stage 2	76.0	0.00
w/o Stage 1	74.3	0.09
w/o Reasoning	76.6	0.09

*Table 6.* Ablation study of InfiGUIAgent components. Scores for the full pipeline are taken from the main paper results for

*Table 5.* Performances on AndroidWorld. Results marked in **bold** represent the best performance

AndroidWorld. Table 5 compares the success rates of *InfiGUIAgent* with open-source models on AndroidWorld (Rawles et al., 2024). InfiGUIAgent-2B achieves an overall success rate of 0.09, outperforming open-source models of similar size, such as ShowUI-2B (Lin et al., 2024) (0.07), and model with much more parameters such as LLaVa-OV-7B (Li et al., 2024a) (0.00) and Qwen2-VL-72B (Bai et al., 2023b) (0.05).

#### 4.3. Ablation Studies

We performed ablation studies (Table 6) to assess the contributions of InfiGUIAgent's main components. The full InfiGUIAgent-2B model achieves 76.8% accuracy on ScreenSpot and 0.13 SR on AndroidWorld.

Removing Stage 2 training ("w/o Stage 2") causes the AndroidWorld SR to drop to 0.00, highlighting that the advanced reasoning skills from this stage are crucial for complex, multi-step tasks. Conversely, omitting the dedicated Stage 1 training ("w/o Stage 1") degrades ScreenSpot accuracy to 74.3%. While this underscores the importance

of Stage 1 for foundational GUI understanding and grounding, the model's residual performance on ScreenSpot in this "w/o Stage 1" setting is likely supported by the Stage 2 training data, which includes a "Stage 1-aligned" component composed of data thematically similar to that used in Stage 1. The "w/o Stage 1" configuration still achieves 0.09 SR on AndroidWorld, suggesting that advanced reasoning can offer some utility even with a less robust dedicated perceptual foundation, but is outperformed by the full model.

The "w/o Reasoning" variant, representing a simplified Stage 2 without the specific hierarchical and expectationreflection structures, yields 76.6% on ScreenSpot and 0.09 SR on AndroidWorld. This indicates that while any learning in Stage 2 is beneficial for complex tasks compared to removing Stage 2 entirely, the structured reasoning patterns contribute to the full model's superior performance. These results affirm the distinct and vital roles of both Stage 1 and Stage 2.

#### 4.4. Qualitative Analysis

A qualitative analysis of the models' performance on the stopwatch operation task reveals significant disparities in their capabilities. The baseline model (Qwen2-VL-2B) exhibits behavior consistent with random exploration, demonstrating a fundamental lack of understanding regarding the screen state or the overarching task goal. As detailed in Table 7, its sequence of three taps appears arbitrary, failing to engage with crucial UI elements such as the "Stopwatch" tab or the "Start" button. This deficiency in identifying and interacting with relevant components results in an inability to progress, leaving the task uncompleted, and suggests an absence of visual grounding and goal-oriented reasoning.

In stark contrast, InfiGUIAgent demonstrates a structured and effective multi-step approach, visually and descriptively detailed in Table 7. For instance, in its initial step, the agent correctly identifies the need to navigate to the Stopwatch tab and executes the tap, expecting the UI to transition. Subsequent steps involve confirming the UI state, identifying the next sub-goal, locating the relevant UI element for that sub-goal, and performing the action with a clear expectation. Each step showcases its ability to process the visual information from the screen (represented by the included images for each step) and make an informed decision.

This side-by-side qualitative assessment underscores InfiGUIAgent's markedly superior capabilities in robust visual grounding, effective hierarchical task decomposition, and strong goal-oriented reasoning. The agent's explicit articulation of its internal reasoning, planned actions, and expected outcomes at each juncture, supported by the visual context of each step, highlights a sophisticated and transparent understanding of the GUI interaction process, ultimately leading to successful task execution where the baseline model fails. See Table 7 in Appendix for more details.

## 5. Conclusion

In this work, we propose *InfiGUIAgent*, a novel MLLMbased GUI Agents. By constructing comprehensive training datasets with two-stage supervised fine-tuning, we enhance the model's ability to understand, reason, and interact with GUIs. Our evaluation, conducted using raw screenshots without relying on additional GUI metadata, demonstrates the model's applicability to real-world scenarios. Experimental results show that our model performs well on GUI tasks and surpass several open-source baselines. We markedly enhance the model's core abilities in visual GUI understanding, multi-step task reasoning, and precise interactive control.

#### **Impact Statement**

This paper presents work whose goal is to advance the field of GUI Agents. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

### References

- Alayrac, J.-B., Donahue, J., Luc, P., Miech, A., Barr, I., Hasson, Y., Lenc, K., Mensch, A., Millican, K., Reynolds, M., et al. Flamingo: a visual language model for fewshot learning. *Advances in neural information processing* systems, 35:23716–23736, 2022.
- Baechler, G., Sunkara, S., Wang, M., Zubach, F., Mansoor, H., Etter, V., Cărbune, V., Lin, J., Chen, J., and Sharma, A. Screenai: A vision-language model for ui and infographics understanding. *arXiv preprint arXiv:2402.04615*, 2024.
- Bai, C., Zang, X., Xu, Y., Sunkara, S., Rastogi, A., Chen, J., and y Arcas, B. A. Uibert: Learning generic multimodal representations for ui understanding. *arXiv preprint arXiv:2107.13731*, 2021.
- Bai, J., Bai, S., Chu, Y., Cui, Z., Dang, K., Deng, X., Fan, Y., Ge, W., Han, Y., Huang, F., Hui, B., Ji, L., Li, M., Lin, J., Lin, R., Liu, D., Liu, G., Lu, C., Lu, K., Ma, J., Men, R., Ren, X., Ren, X., Tan, C., Tan, S., Tu, J., Wang, P., Wang, S., Wang, W., Wu, S., Xu, B., Xu, J., Yang, A., Yang, H., Yang, J., Yang, J., Yang, S., Yao, Y., Yu, B., Bowen, Y., Yuan, H., Yuan, Z., Zhang, J., Zhang, X., Zhang, Y., Zhang, Z., Zhou, C., Zhou, J., Zhou, X., and Zhu, T. Qwen technical report. *ArXiv*, 2023a. URL https://doi.org/10.48550/arXiv.2309.16609.
- Bai, J., Bai, S., Yang, S., Wang, S., Tan, S., Wang, P., Lin, J., Zhou, C., and Zhou, J. Qwen-vl: A frontier large vision-language model with versatile abilities. ArXiv, 2023b. URL https://doi.org/10. 48550/arXiv.2308.12966.
- Bai, J., Bai, S., Yang, S., Wang, S., Tan, S., Wang, P., Lin, J., Zhou, C., and Zhou, J. Qwen-vl: A frontier large visionlanguage model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023c.
- Chai, Y., Huang, S., Niu, Y., Xiao, H., Liu, L., Zhang, D., Gao, P., Ren, S., and Li, H. Amex: Android multiannotation expo dataset for mobile gui agents. arXiv preprint arXiv:2407.17490, 2024.
- Chen, W., Cui, J., Hu, J., Qin, Y., Fang, J., Zhao, Y., Wang, C., Liu, J., Chen, G., Huo, Y., Yao, Y., Lin, Y., Liu, Z., and Sun, M. Guicourse: From general vision language models to versatile gui agents. *arXiv preprint arXiv:2406.11317*, 2024.

- Cheng, K., Sun, Q., Chu, Y., Xu, F., Li, Y., Zhang, J., and Wu, Z. Seeclick: Harnessing gui grounding for advanced visual gui agents. arXiv preprint arXiv:2401.10935, 2024.
- Dai, Y., Tang, D., Liu, L., Tan, M., Zhou, C., Wang, J., Feng, Z., Zhang, F., Hu, X., and Shi, S. One model, multiple modalities: A sparsely activated approach for text, sound, image, video and code, 2022. URL https: //arxiv.org/abs/2205.06126.
- Dao, T. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv preprint arXiv:2307.08691*, 2023.
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby, N. An image is worth 16x16 words: Transformers for image recognition at scale. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021. URL https://openreview.net/forum? id=YicbFdNTTy.
- Floridi, L. and Chiriatti, M. Gpt-3: Its nature, scope, limits, and consequences. *Minds and Machines*, 30:681–694, 2020.
- Glaive AI. Glaive function calling dataset. https://huggingface.co/datasets/ glaiveai/glaive-function-calling, 2024. Accessed: 2024-01-08.
- Gou, B., Wang, R., Zheng, B., Xie, Y., Chang, C., Shu, Y., Sun, H., and Su, Y. Navigating the digital world as humans do: Universal visual grounding for gui agents. *arXiv preprint arXiv:2410.05243*, 2024.
- Hong, W., Wang, W., Lv, Q., Xu, J., Yu, W., Ji, J., Wang, Y., Wang, Z., Dong, Y., Ding, M., et al. Cogagent: A visual language model for gui agents. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14281–14290, 2024.
- Hsiao, Y.-C., Zubach, F., Baechler, G., Carbune, V., Lin, J., Wang, M., Sunkara, S., Zhu, Y., and Chen, J. Screenqa: Large-scale question-answer pairs over mobile app screenshots. *arXiv preprint arXiv:2209.08199*, 2022.
- Hu, X., Kuang, K., Sun, J., Yang, H., and Wu, F. Leveraging print debugging to improve code generation in large language models, 2024a. URL https://arxiv.org/ abs/2401.05319.
- Hu, X., Xiong, T., Yi, B., Wei, Z., Xiao, R., Chen, Y., Ye, J., Tao, M., Zhou, X., Zhao, Z., Li, Y., Xu, S., Wang, S., Xu, X., Qiao, S., Kuang, K., Zeng, T., Wang, L.,

Li, J., Jiang, Y. E., Zhou, W., Wang, G., Yin, K., Zhao, Z., Yang, H., Wu, F., Zhang, S., and Wu, F. Os agents: A survey on mllm-based agents for general computing devices use. *Preprints*, December 2024b. doi: 10.20944/ preprints202412.2294.v1. URL https://doi.org/ 10.20944/preprints202412.2294.v1.

- Hu, X., Zhao, Z., Wei, S., Chai, Z., Ma, Q., Wang, G., Wang, X., Su, J., Xu, J., Zhu, M., Cheng, Y., Yuan, J., Li, J., Kuang, K., Yang, Y., Yang, H., and Wu, F. Infiagentdabench: Evaluating agents on data analysis tasks. *arXiv* preprint arXiv:2401.05507, 2024c.
- Huang, J. and Chang, K. C.-C. Towards reasoning in large language models: A survey, 2023. URL https:// arxiv.org/abs/2212.10403.
- Huang, X., Liu, W., Chen, X., Wang, X., Wang, H., Lian, D., Wang, Y., Tang, R., and Chen, E. Understanding the planning of llm agents: A survey. *arXiv preprint arXiv:2402.02716*, 2024.
- Jiang, Y., Schoop, E., Swearngin, A., and Nichols, J. Iluvui: Instruction-tuned language-vision modeling of uis from machine conversations. *arXiv preprint arXiv:2310.04869*, 2023.
- Kapoor, R., Butala, Y. P., Russak, M., Koh, J. Y., Kamble, K., Alshikh, W., and Salakutdinov, R. Omniact: A dataset and benchmark for enabling multimodal generalist autonomous agents for desktop and web. arXiv preprint arXiv:2402.17553, 2024.
- Li, B., Zhang, Y., Guo, D., Zhang, R., Li, F., Zhang, H., Zhang, K., Zhang, P., Li, Y., Liu, Z., and Li, C. Llavaonevision: Easy visual task transfer, 2024a. URL https: //arxiv.org/abs/2408.03326.
- Li, B., Zhang, Y., Guo, D., Zhang, R., Li, F., Zhang, H., Zhang, K., Zhang, P., Li, Y., Liu, Z., and Li, C. Llavaonevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03329*, 2024b.
- Li, D., Liu, Y., Wu, H., Wang, Y., Shen, Z., Qu, B., Niu, X., Zhou, F., Huang, C., Li, Y., Zhu, C., Ren, X., Li, C., Ye, Y., Zhang, L., Yan, H., Wang, G., Chen, B., and Li, J. Aria: An open multimodal native mixture-ofexperts model, 2024c. URL https://arxiv.org/ abs/2410.05993.
- Li, J., Li, D., Savarese, S., and Hoi, S. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference* on machine learning, pp. 19730–19742. PMLR, 2023.
- Li, L., Geng, S., Li, Z., He, Y., Yu, H., Hua, Z., Ning, G., Wang, S., Xie, T., and Yang, H. Infibench: Evaluating the question-answering capabilities of code large language models. *arXiv preprint arXiv:2404.07940*, 2024d.

- Li, Y., He, J., Zhou, X., Zhang, Y., and Baldridge, J. Mapping natural language instructions to mobile ui action sequences. *arXiv preprint arXiv:2005.03776*, 2020a.
- Li, Y., Li, Gangaand He, L., Zheng, J., Li, H., and Guan, Z. Widget captioning: Generating natural language description for mobile user interface elements. *arXiv preprint arXiv:2010.04295*, 2020b.
- Lin, K. Q., Li, L., Gao, D., Yang, Z., Bai, Z., Lei, W., Wang, L., and Shou, M. Z. Showui: One vision-languageaction model for generalist gui agent. In *NeurIPS 2024 Workshop on Open-World Agents*, 2024.
- Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tuning. Advances in neural information processing systems, 36:34892–34916, 2023.
- Liu, H., You, Q., Wang, Y., Han, X., Zhai, B., Liu, Y., Chen, W., Jian, Y., Tao, Y., Yuan, J., He, R., and Yang, H. Infimm: Advancing multimodal understanding with an open-sourced visual language model. In *Annual Meeting* of the Association for Computational Linguistics, 2024.
- Liu, Z., Mao, H., Wu, C.-Y., Feichtenhofer, C., Darrell, T., and Xie, S. A convnet for the 2020s. In *Proceedings of* the IEEE/CVF conference on computer vision and pattern recognition, pp. 11976–11986, 2022.
- MDeitke, M., Clark, C., Lee, S., Tripathi, R., Yang, Y., Park, J. S., Salehi, M., Muennighoff, N., Lo, K., Soldaini, L., Lu, J., Anderson, T., Bramsom, E., Ehsani, K., Ngo, H., Chen, Y., Patel, A., Yatskar, M., Callison-Burch, C., Head, A., Hendrix, R., Bastani, F., van der Bilt, E., Lambert, N., Chou, Y., Chheda, A., Sparks, J., Skjonsberg, S., Schmitz, M., Sarnat, A., Bischoff, B., Walsh, P., Newell, C., Wolters, P., Gupta, Tanmay sna Zeng, K.-H., Borchardt, J., Groeneveld, D., Nam, C., Lebrecht, S., Wittlif, C., Schoenick, C., Michel, O., Krishna, R., Weihs, L., Smith, N. A., Hajishirzi, H., Girshick, R., Farhadi, A., and Kembhavi, A. Molmo and pixmo: Open weights and open data for state-of-the-art vision-language models. *arXiv preprint arXiv:2409.17146*, 2024.
- Nong, S., Zhu, J., Wu, R., Jin, J., Shan, S., Huang, X., and Xu, W. Mobileflow: A multimodal llm for mobile gui agent. arXiv preprint arXiv:2407.04346, 2024.
- OpenAI. Gpt-4v(ision) system card, 2023. URL https://cdn.openai.com/papers/GPTV\_ System\_Card.pdf.
- OpenAI. Gpt-4o, 2024. URL https://openai.com/ index/hello-gpt-4o/. Accessed: 2025-01-03.
- Qi, Z., Liu, X., Iong, I. L., Lai, H., Sun, X., Yang, X., Sun, J., Yang, Y., Yao, S., Zhang, T., et al. Webrl: Training

llm web agents via self-evolving online curriculum reinforcement learning. *arXiv preprint arXiv:2411.02337*, 2024.

- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PMLR, 2021.
- Rajbhandari, S., Rasley, J., Ruwase, O., and He, Y. Zero: Memory optimizations toward training trillion parameter models, 2020. URL https://arxiv.org/abs/ 1910.02054.
- Rawles, C., Clinckemaillie, S., Chang, Y., Waltz, J., Lau, G., Fair, M., Li, A., Bishop, W., Li, W., Campbell-Ajala, F., et al. Androidworld: A dynamic benchmarking environment for autonomous agents. *arXiv preprint arXiv:2405.14573*, 2024.
- Shinn, N., Cassano, F., Berman, E., Gopinath, A., Narasimhan, K., and Yao, S. Reflexion: Language agents with verbal reinforcement learning, 2023. URL https://arxiv.org/abs/2303.11366.
- Sunkara, S., Wang, M., Liu, L., Baechler, G., Hsiao, Y.-C., Sharma, A., Stout, J., et al. Towards better semantic understanding of mobile interfaces. *arXiv preprint arXiv:2210.02663*, 2022.
- Team, G., Georgiev, P., Lei, V. I., Burnell, R., Bai, L., Gulati, A., Tanzer, G., Vincent, D., Pan, Z., Wang, S., et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. arXiv preprint arXiv:2403.05530, 2024.
- Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., and Lample, G. Llama: Open and efficient foundation language models. *ArXiv*, 2023. URL https: //doi.org/10.48550/arXiv.2302.13971.
- Vu, M. D., Wang, H., Chen, J., Li, Z., Zhao, S., Xing, Z., and Chen, C. Gptvoicetasker: Advancing multi-step mobile task efficiency through dynamic interface exploration and learning. In *Proceedings of the 37th Annual ACM Symposium on User Interface Software and Technology*, pp. 1–17, 2024.
- Wang, B., Li, G., Zhou, X., Chen, Z., Grossman, T., and Li, Y. Screen2words: Automatic mobile ui summarization with multimodal learning. arXiv preprint arXiv:2108.03353, 2021.

- Wang, P., Bai, S., Tan, S., Wang, S., Fan, Z., Bai, J., Chen, K., Liu, X., Wang, J., Ge, W., et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. arXiv preprint arXiv:2409.12191, 2024.
- Wang, W., Lv, Q., Yu, W., Hong, W., Qi, J., Wang, Y., Ji, J., Yang, Z., Zhao, L., Song, X., et al. Cogvlm: Visual expert for pretrained language models. *arXiv preprint arXiv:2311.03079*, 2023.
- Wen, H., Li, Y., Liu, G., Zhao, S., Yu, T., Li, T. J.-J., Jiang, S., Liu, Y., Zhang, Y., and Liu, Y. Autodroid: Llm-powered task automation in android. *arXiv preprint arXiv:2308.15272*, 2023.
- Xiao, C., Hu, X., Liu, Z., Tu, C., and Sun, M. Lawformer: A pre-trained language model for chinese legal long documents. AI Open, 2:79–84, 2021.
- Xie, T., Zhang, D., Chen, J., Li, X., Zhao, S., Cao, R., Hua, T. J., Cheng, Z., Shin, D., Lei, F., et al. Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments. *arXiv preprint arXiv:2404.07972*, 2024.
- Yang, J., Zhang, H., Li, F., Zou, X., Li, C., and Gao, J. Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v. arXiv preprint arXiv:2310.11441, 2023a.
- Yang, J., Zhang, H., Li, F., Zou, X., Li, C., and Gao, J. Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v, 2023b. URL https://arxiv. org/abs/2310.11441.
- Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., and Cao, Y. React: Synergizing reasoning and acting in language models, 2023. URL https://arxiv. org/abs/2210.03629.
- Yin, D., Brahman, F., Ravichander, A., Chandu, K., Chang, K.-W., Choi, Y., and Lin, B. Y. Agent lumos: Unified and modular training for open-source language agents. arXiv preprint arXiv:2311.05657, 2023.
- You, K., Zhang, H., Schoop, E., Weers, F., Swearngin, A., Nichols, J., Yang, Y., and Gan, Z. Ferret-ui: Grounded mobile ui understanding with multimodal llms. In *European Conference on Computer Vision*, pp. 240–255. Springer, 2025.
- Yu, X., Peng, B., Vajipey, V., Cheng, H., Galley, M., Gao, J., and Yu, Z. Exact: Teaching ai agents to explore with reflective-mcts and exploratory learning. *arXiv preprint arXiv:2410.02052*, 2024.
- Zhang, C., Yang, Z., Liu, J., Han, Y., Chen, X., Huang, Z., Fu, B., and Yu, G. Appagent: Multimodal agents

as smartphone users. *arXiv preprint arXiv:2312.13771*, 2023.

- Zhang, J., Wu, J., Teng, Y., Liao, M., Xu, N., Xiao, X., Wei, Z., and Tang, D. Android in the zoo: Chain-of-actionthought for gui agents. *arXiv preprint arXiv:2403.02713*, 2024a.
- Zhang, Y., Mao, S., Ge, T., Wang, X., de Wynter, A., Xia, Y., Wu, W., Song, T., Lan, M., and Wei, F. Llm as a mastermind: A survey of strategic reasoning with large language models, 2024b. URL https://arxiv.org/ abs/2404.01230.
- Zhang, Z. and Zhang, A. You only look at screens: Multimodal chain-of-action agents. *arXiv preprint arXiv:2309.11436*, 2023.

## A. Limitations

While our work demonstrates promising advancements in GUI task automation through the proposed two-stage training framework, several limitations remain. First, due to our focus on enabling efficient deployment on edge devices, we primarily explored small-scale models (e.g., Qwen2-VL-2B). While this approach ensures practicality, larger models may offer further performance improvements, particularly in handling more complex tasks. Second, although hierarchical reasoning and expectation-reflection reasoning enhance the agent's ability to decompose tasks and adapt dynamically, real-world GUI interactions often involve unforeseen complexities, such as error recovery, cross-application coordination, or dynamically changing interfaces. Enhancing the agent's robustness to such challenges remains an open research direction.

## **B.** Prompt Templates

This section details the various prompt templates employed for trajectory and grounding tasks within our study. These prompts are designed to guide the AI agent in understanding its role, the environment, and the specific requirements of each task.

## **B.1. Prompt Templates for Grounding Tasks**

For grounding tasks, which require the agent to identify specific elements or locations on the UI, we use tailored prompt templates. The choice of template depends on whether the output is a single point or a bounding box.

## B.1.1. POINT OUTPUT

When the grounding task requires the output of specific coordinates (a point) related to a given instruction, the following prompt template is used:

## **Prompt Template for Grounding (Point)**

Output the relative coordinates of the icon, widget, or text most closely related to "instruction" in this screenshot, in the format of " $\{"x": x, "y": y\}$ ", where x and y are in the positive directions of horizontal left and vertical down respectively, with the origin at the top left corner, and the range is 0-1000.

Here, {instruction} is a placeholder for the natural language instruction describing the target element (e.g., "the 'Login' button", "the user's profile picture").

## B.1.2. BOUNDING BOX OUTPUT

When the grounding task requires identifying a bounding box for an element related to a given instruction, the following prompt template is employed:

## Prompt Template for Grounding (Bounding Box)

Output the relative coordinates of the icon, widget, or text most closely related to "instruction" in this screenshot, in the format of " $\{x1": x1, y1": y1, x2": x2, y2": y2\}$ ", where x1, y1, x2 and y2 are in the positive directions of horizontal left and vertical down respectively, with the origin at the top left corner, and the range is 0-1000.

Similarly, {instruction} serves as a placeholder for the user's instruction. The coordinates x1, y1 represent the top-left corner of the bounding box, and x2, y2 represent the bottom-right corner.

## **B.2. System Prompt for Trajectory Tasks**

For trajectory-based tasks, which involve the agent performing a sequence of actions to achieve a goal on a user interface, we utilize the following comprehensive system prompt:

#### System Prompt for Trajectory Tasks

You are an AI agent capable of interacting with a user interface through function calling. Your task is to assist users in completing UI-related tasks or answering related questions.

The current UI type is: {ui\_type}

When referring to positions on the UI:

- Use integer relative coordinates ranging from 0 to 1000.
- The origin (0, 0) is at the top-left corner of the screen. The x-axis increases to the right, and the y-axis increases downward.

You are provided with function signatures within  $i_{tools_{\ell}}$   $j_{tools_{\ell}}$  XML tags. You may call one or more functions to assist with the UI-related tasks. Don't make assumptions about what values to plug into functions.  $i_{tools_{\ell}}$ 

- {action\_space}
- ;/tools¿

For each function call return a json object with function name and arguments within ¡tool\_call¿ ¡/tool\_call¿ tags with the following schema:

itool\_call;

"name": ¡function-name¿, "arguments": ¡args-dict¿ ¡/tool\_call¿

In this prompt,  $\{ui\_type\}$  is a placeholder that is replaced with the actual type of the user interface being interacted with. The available actions and their function signatures are provided to the agent within the tools and /tools tags at runtime.

## C. Cases

## C.1. Stage 1: Fundamental Abilities

We demonstrate the fundamental abilities trained in Stage 1 through three cases: GUI Understanding (Figure 3), Grounding (Figure 4), and Question Answering (Figure 5).

## C.2. Stage 2: Native Reasoning

We provide two representative cases to demonstrate the reasoning and interaction process of InfiGUIAgent.

**Reply to a Message** Figure 6 illustrates a step where the agent needs to reply to a specific message in a messaging application. The reasoning process involves identifying the "Start chat" button and grounding the action to initiate the reply process.

**Creating a New Contact** Figure 7 and Figure 8 demonstrate sequential steps for creating a new contact. In the first step (Step K), the agent navigates to the "Contacts" section by reasoning and grounding the action to the corresponding tab. In the following step (Step K+1), the agent initiates the contact creation process by identifying and tapping the "Create new contact" button. These sequential steps highlight the agent's hierarchical reasoning and grounding abilities.

Baseline Output (Qwen2-VL-2B, temp=1.0)	InfiGUIAgent Output (Stage 1 + 2, temp=1.0)
Actions: tap(x=710, y=55) tap(x=500, y=880)	Step 1: Navigate to Stopwatch Tab
tap(x=240, y=240)	8 +
<b>Provided Observation:</b> Issues taps without clear reasoning. Fails to locate Stopwatch tab or start button; never reaches running state. No goal conditioning or hierarchical planning evident.	8/ 8/ 8/ 8/ 8/ 8/ 8/ 8/ 8/ 8/
	<ul> <li>Agent Reasoning (Key Points): Clock app open, tabs visible.</li> <li>Goal: Run stopwatch. Need to:</li> <li>Tap Stopwatch tab.</li> <li>Wait for UI transition.</li> <li>Stopwatch tab is rightmost on top bar.</li> <li>Action: tap (x=800, y=120)</li> <li>Expected Outcome: Stopwatch screen loads with a start button.</li> </ul>
	Step 2: Start the Stopwatch
	Agent Reasoning (Key Points): Now on Stopwatch screen. Timer and Start button visible. Goal: Start stopwatch. Need to tap Start button. Start button is large, circular, centered below timer. Action: tap ( $x=500$ , $y=850$ ) Expected Outcome: Timer starts counting, stopwatch running.
	Step 3: Stop the Stopwatch
	Agent Reasoning (Key Points): Stopwatch is running. Next goal is to stop it. "Stop" button is now visible (often replaces "Start"). Action: tap(x=500, y=850) Expected Outcome: Timer stops counting. "Reset" and "Lap" buttons might appear.
	Step 4: Reset the Stopwatch
	Agent Reasoning (Key Points): Stopwatch is stopped. Final goal is to reset. "Reset" button should be visible. Action: $tap(x=250, y=850)$ Expected Outcome: Timer resets to 00:00.00. "Start" button reappears.

Table 7. Qualitative Comparison for Task: ClockStopWatchRunning (Illustrating a 4-step InfiGUIAgent process with screen states)



Figure 3. Case of GUI Understanding.



Figure 4. Case of Grounding.



Figure 5. Case of Question Answering.



Figure 6. Case of Native Advanced Reasoning. The agent's goal is to reply to a message



Figure 7. Case of Native Advanced Reasoning. The agent's goal is to create a new contact.



Figure 8. Case of Native Advanced Reasoning. The agent's goal is to create a new contact.