



# FERRET-UI ONE: MASTERING UNIVERSAL USER INTERFACE UNDERSTANDING ACROSS PLATFORMS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Building a generalist model for user interface (UI) understanding is challenging due to various foundational issues, such as platform diversity, resolution variation, and data limitation. In this paper, we introduce **Ferret-UI One**, a multimodal large language model (MLLM) designed for universal UI understanding across a wide range of platforms, including iPhone, Android, iPad, Webpage, and AppleTV. Building on the foundation of Ferret-UI, Ferret-UI One introduces three key innovations: support for multiple platform types, high-resolution perception through adaptive scaling, and advanced task training data generation powered by GPT-4o with set-of-mark visual prompting. These advancements enable Ferret-UI One to perform complex, user-centered interactions, making it highly versatile and adaptable for the expanding diversity of platform ecosystems. Extensive empirical experiments on referring, grounding, user-centric advanced tasks (comprising 9 subtasks  $\times$  5 platforms), GUIDE next-action prediction dataset, and GUI-World multi-platform benchmark demonstrate that Ferret-UI One significantly outperforms Ferret-UI, and also shows strong cross-platform transfer capabilities.

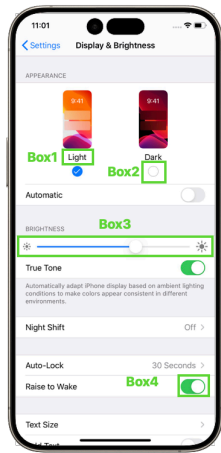
## 1 INTRODUCTION

User interfaces (UIs) are central to human-computer interaction, shaping how users interact with digital systems. The complexity of UIs has evolved with the proliferation of platforms such as smartphones, tablets, web platforms, and smart TVs. Despite this increasing diversity, many current approaches to UI understanding and interaction (Hong et al., 2023; Wang et al., 2024a; Kapoor et al., 2024), particularly those in multi-platform ecosystems, face limitations.

One prominent effort in this space is Ferret-UI (You et al., 2024), which has advanced the field of referring and grounding UIs. However, though taking an any-resolution approach (Liu et al., 2024a; Zhang et al., 2024b), Ferret-UI is constrained by a fixed grounding resolution (*i.e.*,  $336 \times 672$  and  $672 \times 336$ ), and focuses on single-type platforms (*i.e.*, mobile devices including iPhone and Android), limiting its applicability in the context of today’s highly diverse platform landscape. For example, as illustrated in Figure 1, one notable difference among these four exemplified platforms is resolution, the native resolution of iPhone differs from that of iPad, Web UI and also AppleTV, directly applying Ferret-UI across these platforms presents significant challenges. Another major challenge is the lack of platform-specific, high-quality data, given different platforms. Though Ferret-UI’s approach to task training data generation can be extended to these platforms, it primarily relies on text-based GPT-4 prompting, where bounding boxes are represented in a purely textual format. This absence of visual input and spatial relationships between UI elements diminishes the quality of training data, which in turn limits the performance and effectiveness of the resulting model.

To address these limitations, we introduce **Ferret-UI One**, a multimodal large language model (MLLM) designed to understand diverse UI screens and respond to user intent through single-step interactions across multiple platforms. Building on the foundation of Ferret-UI (You et al., 2024), Ferret-UI One significantly enhances UI perception and user interaction capabilities via three key advancements: (*i*) multi-platform support, (*ii*) dynamic high-resolution image encoding, and (*iii*) high-quality multimodal training data generation.

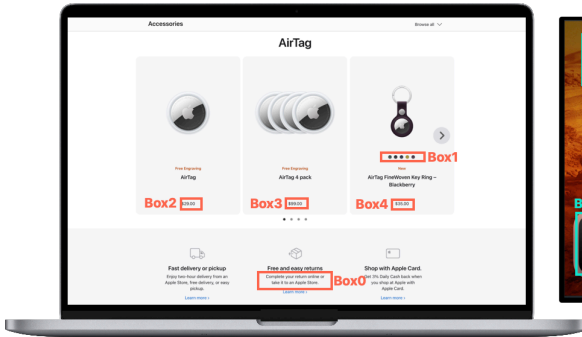
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107



- What kind of UI widget is in [Box0]?
- This UI widget is a toggle.
- What is the current state of the appearance toggle and how to switch it?
- The current theme is set to Light Mode [Box1]. To switch to Dark Mode, you would tap on the radio button [Box2] next to "Dark", which is currently unchecked.
- How would you interact with the brightness slider to decrease the screen brightness?
- The current brightness level is slightly above mid-level. To decrease the brightness, drag the slider [Box3] towards the left, where the lower brightness icon is located.



- What does the text in [Box0] indicate?
- The text that reads "10-DAY FORECAST" indicates the list below shows the weather forecast of the next 10 days.
- What is the current temperature in Cupertino, and what is the weather condition?
- The current temperature in Cupertino is 91°F [Box1], and the weather condition is "Mostly Sunny" [Box2].
- How to view the weather conditions of later hours?
- You can scroll horizontally on the hourly forecast row [Box3] to see weather conditions for later hours.



- What does the text in [Box0] say?
- The text says "Complete your return online or take it to an Apple Store".
- How to explore different color options for the AirTag FineWoven Key Ring?
- Tap or click on the dots [Box1] below the image of the AirTag FineWoven Key Ring.
- What are the prices of the three visible products on the screen?
- The prices are AirTag for \$29.00 [Box2], AirTag 4-pack for \$99.00 [Box3], and AirTag FineWoven Key Ring - Blackberry for \$35.00 [Box4].



- What is the [Box0] icon?
- This icon is for "App Store" on Apple's devices.
- What does the large "For All Mankind" banner suggest about the context of this screen?
- The large banner suggests that "For All Mankind" [Box1] is either a featured show or content available on the Apple TV app. The banner acts as promotional content, encouraging users to watch or explore the show.
- Where do I open the Apple TV app?
- To open the Apple TV app, tap or select the black tile with the Apple TV logo [Box2].

Figure 1: Real examples of a single Ferret-UI One model interacting with four different platforms (iPhone, iPad, Webpage, and AppleTV) for UI understanding.

108 First, Ferret-UI One extends its compatibility beyond mobile platforms (iPhone and Android), incor-  
 109 porating additional platforms like tablets, webpages, and smart TVs. Figure 1 illustrates visual ex-  
 110 amples of Ferret-UI One interacting with users across four typical screen types. This multi-platform  
 111 support enables broader applicability and allows the system to seamlessly scale across a variety of  
 112 user environments.

113 Second, Ferret-UI One supports high-resolution image encoding via the any-resolution method (Liu  
 114 et al., 2024a). However, going beyond that, we introduce an enhanced *adaptive gridding* approach  
 115 to maintain perception capabilities at the original resolution of the UI screenshot, ensuring more  
 116 accurate recognition of visual elements. By leveraging human-collected bounding box annotations,  
 117 we enhance referring and grounding precision, which allows for more detailed understanding of UI  
 118 components and their relationships.

119 Third, Ferret-UI One leverages high-quality training data for both elementary and advanced tasks.  
 120 For elementary tasks, we convert simple referring and grounding data into conversations, allowing  
 121 the model to develop a fundamental understanding of diverse UI screens. For advanced tasks, which  
 122 focus on user-centered, free-form conversations, we replace the text-based GPT-4 prompting (where  
 123 bounding boxes are described only in text) with GPT-4o using set-of-mark visual prompting (Yang  
 124 et al., 2023a) for training data generation. This approach enhances spatial understanding of UI  
 125 elements, resulting in higher-quality training data. Additionally, unlike previous methods that use  
 126 straightforward instructions such as “click on [bbox location]”, Ferret-UI One performs single-step  
 127 user-centered interactions. For example, when given a command like “please confirm submission”,  
 128 the system understands and executes the intended action, rather than simply following mechanical  
 129 click instructions. Overall, our contributions are summarized as follows.

- 130 • We present Ferret-UI One, a multimodal LLM that sets itself apart from previous efforts by  
 131 supporting a broader range of platforms, including iPhone, Android, iPad, Webpage, and Ap-  
 132 pleTV. We upgrade Ferret-UI across multiple fronts, including better instruction-tuning data for  
 133 model training, high-resolution image encoding for enhanced performance, and new referring and  
 134 grounding benchmarks tailored for different platforms.
- 135 • We demonstrate that Ferret-UI One advances the UI referring and grounding performance on  
 136 different platforms. On three categories of tasks (referring, grounding, and user-centric advanced  
 137 tasks, comprising 9 subtasks  $\times$  5 platforms), Ferret-UI One outperforms Ferret-UI, and also shows  
 138 competitive performance compared to GPT-4o. Besides, Ferret-UI One also exhibits strong trans-  
 139 fer capabilities across platforms. Finally, Ferret-UI One achieved strong performance on recent  
 140 benchmarks like GUIDE (Chawla et al., 2024) and GUI-World (Chen et al., 2024).

## 141 2 RELATED WORK

142 UI agents have garnered significant attention in recent research, particularly in multimodal models  
 143 that seek to automate complex UI tasks across diverse platforms. Many works have advanced the  
 144 field by tackling specific challenges related to single-platform and multi-platform UI understanding,  
 145 interaction, and automation.

146 **Single-Platform UI Agents.** Single-platform UI agents focus on automating tasks on a specific de-  
 147 vice ecosystem, such as Android, iOS, desktop environments, or webpages. AndroidEnv (Toyama  
 148 et al., 2021) establishes a platform for reinforcement learning in Android environments, allowing  
 149 agents to learn and interact within mobile applications. Similarly, AssistGUI (Gao et al., 2023)  
 150 focuses on automating desktop GUI tasks, while UFO (Zhang et al., 2024a) presents an agent for  
 151 interacting with Windows OS. On the mobile side, AppAgent (Yang et al., 2023b) is specialized  
 152 for smartphone UI interactions. For web-based agents, systems like WebArena (Zhou et al., 2023),  
 153 LASER (Ma et al., 2023) and WebShop (Yao et al., 2022) explore agents that navigate and perform  
 154 tasks within web environments. In contrast, OS-Copilot (Wu et al., 2024) explores the more com-  
 155 plex computer OS interaction. These efforts have significantly improved task-specific automation,  
 156 although their single-platform nature limits cross-platform flexibility.

157 **Multi-Platform UI Agents.** Multi-platform UI agents have emerged to address the growing com-  
 158 plexity of device ecosystems, supporting a variety of devices and platforms, including mobile, web,  
 159 and desktop environments. Recent works like OmniACT (Kapoor et al., 2024) support both desktop  
 160 and web interfaces, while CogAgent (Hong et al., 2023) supports UI navigation on both PC web-  
 161 pages and Android devices. Furthermore, Mind2Web (Deng et al., 2024) and Mobile-Agent (Wang

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215

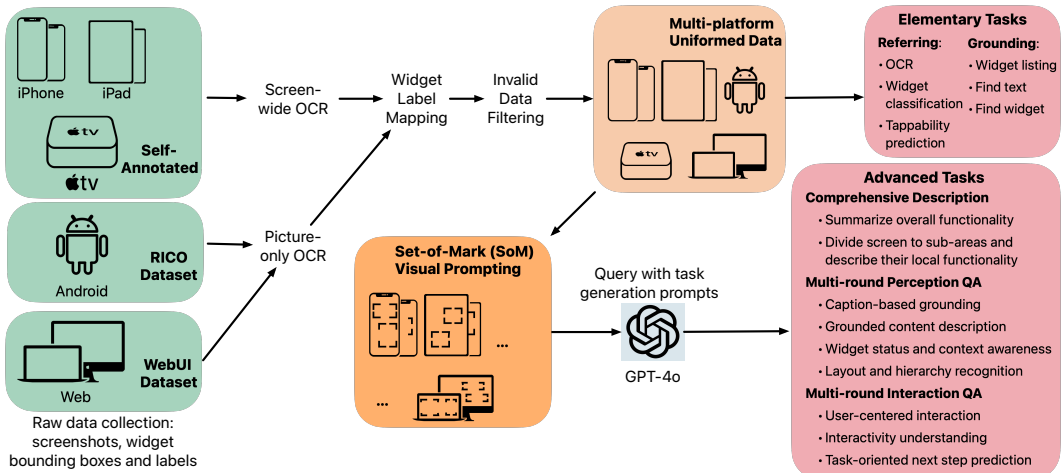


Figure 2: Illustration of the Core-set data generation pipeline.

et al., 2024a) are also notable research works that enable agents to operate seamlessly across different platforms. Ferret-UI (You et al., 2024) focuses on mobile UI understanding for both Android and iPhone screenshots using multimodal LLMs, with a focus on referring and grounding capabilities. These agents aim to perform more complex, user-intent-based interactions across multiple device types, paving the way for truly generalist multimodal agents.

**UI-Agent Benchmarks.** The evaluation of UI agents requires specialized benchmarks to test various aspects of UI interaction, including task execution, navigation, and interaction understanding. Rico (Deka et al., 2017) remains a foundational dataset for mobile app interaction, while Mobile-Env (Zhang et al., 2023), AndroidWorld (Rawles et al., 2024a) and Android in-the-Wild (Rawles et al., 2024b) provide benchmarks for mobile device control. OSWorld (Xie et al., 2024) takes a broader approach by providing benchmarks for agents in real computer environments, including Ubuntu, MacOS, and Windows. Web-based interaction benchmarks include WebSRC (Chen et al., 2021) and WebCanvas (Pan et al., 2024), which focus on structural reading comprehension and task execution in web environments. More recently, MobileAgentBench (Wang et al., 2024b) and VisualWebBench (Liu et al., 2024b) have introduced taxonomies designed to evaluate the performance of multimodal agents across both mobile and web interfaces. VisualAgentBench (Liu et al., 2024c) expands this with a focus on multimodal LLMs as visual foundation agents. These benchmarks contribute to a growing need for unified, cross-platform evaluations that can assess UI agents’ adaptability, precision, and efficiency.

Compared to the aforementioned works, Ferret-UI One is the first to target universal UI understanding across diverse platforms, including smartphones, tablets, web interfaces, and smart TVs. It focuses on foundational capabilities like fine-grained referring, grounding, and reasoning, aiming to create a generalist agent for versatile UI navigation.

### 3 FERRET-UI ONE

In this section, we first describe how we curate our training datasets from the raw data annotations (Section 3.1) and then describe the model architecture (Section 3.2).

#### 3.1 DATASET CONSTRUCTION

We construct our own dataset in order to train a strong multi-platform UI understanding model. A flow diagram of our complete dataset generation pipeline is shown in Figure 2.

**Raw Annotation Collection.** The primary dataset used for training Ferret-UI One is a combination of data from various platform types, including iPhone, Android, iPad, Webpage, and AppleTV. The data collection process varies depending on the platform type:

- **iPhone, iPad, and AppleTV:** We use human-collected iPhone, iPad and AppleTV data under diverse usage scenarios and human-annotated widget bounding box coordinates and labels. To

Table 1: A summary of datasets across various platforms used to train Ferret-UI One. Resolution statistics are shown in .

Training Data	Platform	#Images (k)	Task Types
Core-set	iPhone	112	Elementary Tasks (Referring, Grounding), Advanced Tasks (Comprehensive Description, Multi-Round Perception QA and Interaction QA)
	Android	61	
	iPad	19	
	Web	321	
	AppleTV	16	
GUIDE	Web	51	Next Action Prediction
GroundUI-18k	Web	18	Simple Interaction
Spotlight	Android	66	Screen2Word, Widget Caption, Taperception

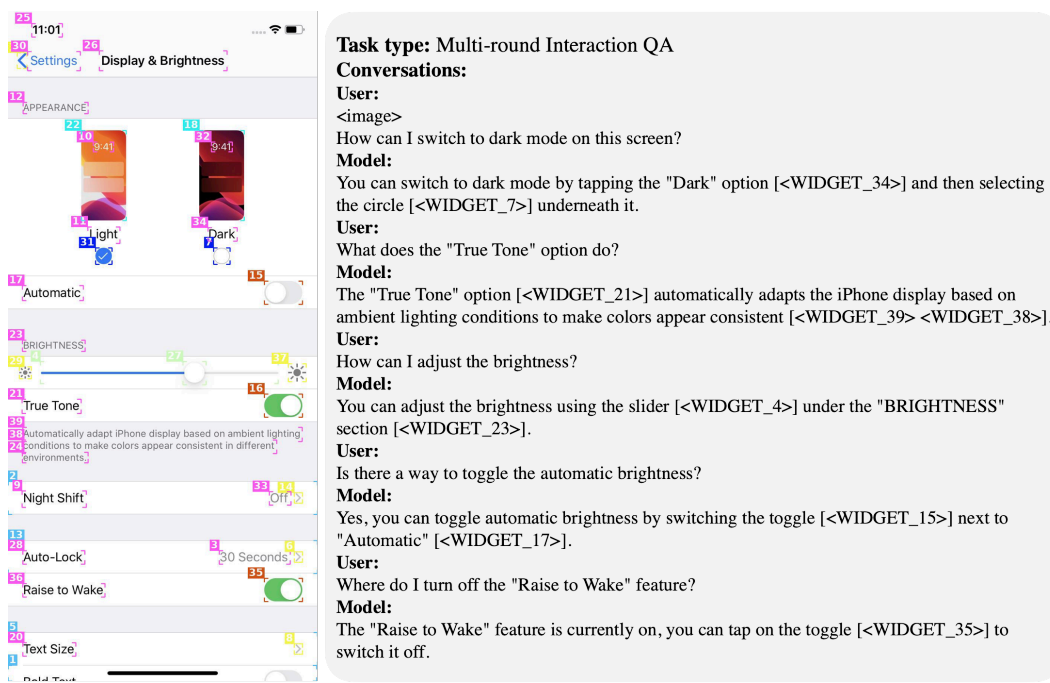


Figure 3: Example of set-of-mark visual prompting (Yang et al., 2023a) (left) and one of its generated advanced task training example (right).

save annotation costs, we do not collect text annotations; instead, text bounding boxes are replaced by screen-wide OCR-detected text and bounding boxes using an OCR confidence threshold of 0.5.

- **Webpage:** The web data is derived from the WebUI dataset (Wu et al., 2023). Bounding boxes of all types of UI widgets and text annotations for non-picture widgets are directly parsed from the source HTML view hierarchy tree, providing high-quality annotations. For picture widgets, we further use OCR to detect texts contained in the pictures.
- **Android:** The Android data for screenshots, bounding boxes and text annotations is transformed from the RICO dataset (Deka et al., 2017). Similar to the WebUI dataset, we also perform picture-only OCR to complete the missing text annotations in picture widgets.

For all the collected data, we perform data filtering including: (i) filter out or narrow down out-of-bound bounding boxes and remove empty screenshots with no remaining bounding boxes after box filtering; (ii) since we do not intend to add multilingual support for the Ferret-UI One model, screenshots with more than 5% non-ASCII characters in the text annotations are also removed.

Despite the different types of label spaces from various sources, we filter out bounding boxes associated with less relevant labels (e.g., UI types) and uniformly map the remaining labels to a common

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323

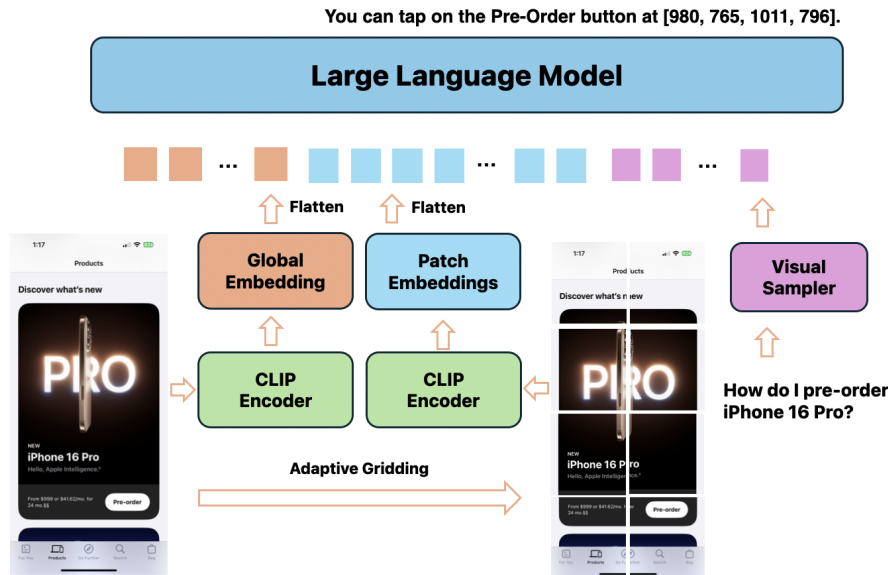


Figure 4: Overview of the Ferret-UI One model architecture, which allows for seamless UI understanding and user-centered single-step interactions with high-resolution support.

label space containing 13 classes: ‘Checkbox’, ‘Button’, ‘Container’, ‘Dialog’, ‘Icon’, ‘PageControl’, ‘Picture’, ‘SegmentedControl’, ‘Slider’, ‘TabBar’, ‘Text’, ‘TextField’, and ‘Toggle’, and obtain a multi-platform uniformed dataset with raw UI widget annotations. We provide the original label statistics and their converted labels in appendix E.

We name the above-collected screenshot dataset the Core-set, which we will use to construct the elementary and advanced task data. Besides, we also employ third-party training datasets to enrich our data source and avoid overfitting our predefined tasks. A complete statistics of the training dataset of Ferret-UI One is summarized in Table 1, indicating that the dataset distribution is very unbalanced across different platforms. In particular, the number of iPad and AppleTV screenshots is significantly smaller than that of other platforms. To address this, we (i) assign different loss weights to different platforms during training, and (ii) generate all three types of advanced tasks for each example of iPad and AppleTV platforms and generate only 1 type of advanced task for each example of other platforms.

Compared to the training dataset of Ferret-UI, which relies on model-detected bounding boxes, Ferret-UI One’s training dataset predominantly utilizes either human-collected annotations or bounding boxes directly parsed from the source HTML, resulting in a significant improvement in annotation quality, evidenced later in our quantitative evaluation in Section 4.2.

**Task Data Generation.** For task data generation, we follow the paradigm of Ferret-UI data construction, which includes both elementary tasks and advanced tasks.

Elementary tasks (Figure 2) consist of 3 referring tasks and 3 grounding tasks. Specifically, referring tasks include (i) *OCR*: recognizing the text given a text bounding box, (ii) *widget classification*: predict the UI type of the elements, and (iii) *tappability*: predict whether the selected widget is tappable for interaction; meanwhile, grounding tasks includes (i) *widget listing*: list all the widgets in the screen, (ii) *find text*: find the location of a given text, and (iii) *find widget*: find the widget given the widget description.

For advanced tasks, we prompt GPT-4o with bounding box annotations of a given screenshot and require GPT-4o to generate QA tasks related to the UI widgets in the screenshot. Unlike Ferret-UI, which mainly focuses on spatial descriptions due to the limitation of using textual prompts without image information (i.e., screenshots) for bounding box annotations, Ferret-UI One leverages GPT-4o to generate advanced task data that covers a variety of aspects of UI understanding. This is possible because GPT-4o demonstrates an improved ability to comprehend the spatial relationships

Algorithm 1: Adaptive  $N$ -gridding

---

**Require:** Original resolution:  $w \times h$ , grid size:  $336 \times 336$ , size limit  $N$   
**Ensure:** Optimal gridding size  $N_w$  and  $N_h$  ( $N_w, N_h \in \mathbb{N}^+$ )

- 1:  $N_{w_{\text{best}}}, N_{h_{\text{best}}} \leftarrow 0, \Delta_{\text{best}} \leftarrow \infty, N_{w_0} \leftarrow \frac{w}{336}, N_{h_0} \leftarrow \frac{h}{336}$
- 2: **for**  $N_w = 1$  to  $N$  **do** ▷ Traverse all grid configurations
- 3:   **for**  $N_h = 1$  to  $N - N_w$  **do**
- 4:      $\Delta_{\text{aspect}} \leftarrow \sqrt{\left| \frac{N_w}{N_h} - \frac{N_{w_0}}{N_{h_0}} \right| \left| \frac{N_h}{N_w} - \frac{N_{h_0}}{N_{w_0}} \right|}$  ▷ Get aspect ratio change
- 5:      $\Delta_{\text{pixel}} \leftarrow \frac{|N_w \times N_h - N_{w_0} \times N_{h_0}|}{N_{w_0} \times N_{h_0}}$  ▷ Get relative pixel change for resizing
- 6:     **if**  $\Delta_{\text{best}} > \Delta_{\text{aspect}} \times \Delta_{\text{pixel}}$  **then**
- 7:        $(N_{w_{\text{best}}}, N_{h_{\text{best}}}) \leftarrow (N_w, N_h)$
- 8:        $\Delta_{\text{best}} \leftarrow \Delta_{\text{aspect}} \times \Delta_{\text{pixel}}$
- 9:     **end if**
- 10:   **end for**
- 11: **end for**
- 12: **return**  $(N_{w_{\text{best}}}, N_{h_{\text{best}}})$

---

between UI widgets when provided with the screenshot as input. Specifically, we prompt GPT-4o to generate 3 types of advanced tasks (shown in Figure 2) including: (i) *comprehensive description*: describe global and local functionalities of the screen, (ii) *multi-round perception QA*: multi-round question answering regarding the UI perception capabilities, and (iii) *multi-round interaction QA*: multi-round question answering regarding the single-step and user-centric UI interactions based on the current screen status. More detailed requirements and prompts for GPT-4o when generating advanced tasks are provided in Appendix A.

We empirically find it hard for GPT-4o to find the location of referred UI widgets with the original screenshot as input (*i.e.*, bad grounding capability). To address this, we use Set-of-Mark (SoM) visual prompting (Yang et al., 2023a) when generating multi-round perception and interaction QA training samples. One example of SoM prompting and its generated data sample is shown in Figure 3, where each UI widget is marked with a corner-style bounding box and a unique number tag for easy identification. Furthermore, the same class of UI widgets have the same color for visual prompting to help GPT-4o better differentiate the bounding boxes of spatially close or nested widgets. Please refer to Figure 5 for visual prompts on other platforms, and Appendix C for additional examples of generated data for advanced tasks.

As aforementioned in Table 1, in addition to the tasks generated on the Core-set, we augment training data with additional third-party datasets, including GroundUI-18k (Zheng et al., 2024), GUIDE (Chawla et al., 2024) and Spotlight (Li & Li, 2023).

### 3.2 MODEL ARCHITECTURE

As shown in Figure 4, the model architecture of Ferret-UI One directly builds upon Ferret-UI (You et al., 2024), which uses the Any-Resolution (AnyRes) method (Liu et al., 2024a) to enhance referring and grounding, enabling the encoder to capture diverse image resolutions.

Specifically, the CLIP image encoder first extracts both global (derived from the low-resolution overview image) and local features (corresponding to high-resolution sub-images) from the UI screenshot. Then, these image features are flattened and sent into the LLM. The Visual Sampler identifies and selects the relevant UI regions based on user instructions. The model then outputs grounded descriptions for perception or interaction with the UI elements.

**Adaptive Gridding.** Local image features are extracted by calculating the optimal grid size using our proposed *adaptive N-gridding* mechanism, then resizing and encoding the visual features of each grid. This is a key model innovation compared to Ferret-UI. As shown in Algorithm 1, the optimal gridding size  $N_w$  and  $N_h$  is determined when the gridding and resizing based on  $(N_w, N_h)$  lead to minimal *aspect ratio change* times *the relative pixel number change*, under the constraint  $N_w + N_h \leq N$ , where  $N$  is the *size limit*. With the size limit  $N$ , the total number of grids is upper bounded by  $\lfloor \frac{N^2}{4} \rfloor$ . Compared to AnyRes module that has unbounded cost, the key differences of

Table 2: Results on our constructed benchmarks for elementary and advanced tasks, as well as the GUIDE benchmark (Chawla et al., 2024). Results on elementary and advanced tasks are averaged over all platforms, including iPhone, Android, iPad, Webpage, and AppleTV. Each platform includes 6 elementary tasks and 3 advanced tasks. (†) In tasks that require referring, GPT-4o is equipped with set-of-mark (SoM) prompting by adding a red rectangular box to screenshots for the referred widget. Note that SoM visual prompting is not used for Ferret-UI and Ferret-UI One.

Model	Backbone	Elementary		Advanced		GUIDE Bench	
		Refer	Ground	GPT-4o Score	Multi-IoU	BertScore	IoU
Ferret-UI	Vicuna-13B	64.15	57.22	45.81	18.75	41.15	26.91
Ferret-UI One	Gemma-2B	75.20	78.13	80.25	40.51	83.71	51.13
	Llama3-8B	80.28	<b>82.79</b>	<b>89.73</b>	41.15	<b>91.37</b>	<b>55.78</b>
	Vicuna-13B	<b>81.34</b>	81.31	86.25	<b>41.71</b>	88.81	54.71
GPT-4o	-	56.47	12.14	77.73	7.06	75.31	9.64
GPT-4o + SoM-Prompt†	-	<b>87.91</b>	-	84.33	7.36	-	-

adaptive  $N$ -gridding is it automatically finds the optimal gridding configuration, i.e. least resolution distortion (aspect ratio change and pixel number change) within a predefined inference cost limit that  $N_w \times N_h \leq \lfloor \frac{N^2}{4} \rfloor$ , which is both information-preserving and efficient for local encoding. With adaptive gridding, Ferret-UI One understands UI screens and provide user-centered interactions with an optimal configuration at any resolution given the inference cost limit specified as size limit  $N$ .

## 4 EXPERIMENTS

### 4.1 EXPERIMENT SETUP

**Training Data.** The training datasets are summarized in Table 1, which can be divided into two categories: (i) datasets constructed by our own, which include elementary task data and advanced task data across all platforms as introduced in Section 3.1, and (ii) public datasets including GroundUI-18k (Zheng et al., 2024), a simple user-centered interaction dataset on webpage screenshots, GUIDE (Chawla et al., 2024), a next-action prediction dataset on webpage screenshots, and Spotlight (Li & Li, 2023), an android UI understanding and interaction dataset.

**Model.** Following Ferret-UI (You et al., 2024), Ferret-UI One uses a CLIP ViT-L/14 model as the image encoder; for the LLM backbone, besides Vicuna-13B (Chiang et al., 2023) as used in the original Ferret-UI, we also tried 2 additional LLMs at mobile scales, including Gemma-2B (Team et al., 2024) and Llama3-8B (Dubey et al., 2024). As to dynamic high-resolution image encoding, we set the size limit  $N$  to 8, so that the maximal grid number is 16 for adaptive gridding.

**Evaluation.** At a high level, model evaluation falls into two broad categories: (i) benchmarks we constructed, and (ii) public benchmarks. For our benchmarks, we created a total of 45, including 6 elementary tasks and 3 advanced tasks per platform type, across 5 platforms. For elementary tasks, we follow the evaluation metrics outlined by You et al. (2023). For advanced tasks, we use GPT-4o to score generated answers for a given screenshot and user query, visually prompting GPT-4o with a red rectangular bounding box. Advanced tasks are tested using GPT-4o evaluation score and multi-IoU. The multi-IoU is calculated by first matching predicted bounding boxes with ground truth bounding boxes and then calculate the average IoU of each pair of bounding boxes (IoU = 0 if no match). Furthermore, we conduct next-action prediction test given previous action history on the GUIDE benchmark (Chawla et al., 2024), and evaluate the semantic similarity w.r.t. the reference answer and grounding Intersection-over-Union (IoU). Additionally, we evaluate our model on the recently released GUI-World benchmark (Chen et al., 2024) on the supported platforms following the original GPT-4 evaluation protocol as in Chen et al. (2024), which does not include evaluating the grounding capability for interaction-related UI tasks.

### 4.2 EXPERIMENT RESULTS

**Main results.** Our main results are summarized in Table 2, which shows the comparative performance of different models on our constructed elementary and advanced tasks, as well as the GUIDE



Table 3: Zero-shot performance of Ferret-UI One on the GUI-World benchmark (Chen et al., 2024).

Model	GPT-4 Score			
	iOS	Android	Webpage	Average
MiniGPT4Video (Ataallah et al., 2024)	1.501	1.342	1.521	1.455
VideoChat2 (Li et al., 2024)	2.169	2.119	2.221	2.170
Chat-Univi (Jin et al., 2024)	2.337	2.390	2.349	2.359
GUI-Vid (Chen et al., 2024)	2.773	2.572	2.957	2.767
QWen-VL-MAX (Bai et al., 2023)	2.779	2.309	2.656	2.580
Ferret-UI (You et al., 2024)	2.713	2.791	2.411	2.638
Ferret-UI One	<b>2.881</b>	<b>2.954</b>	<b>3.013</b>	<b>2.948</b>
Gemini-Pro 1.5 (Reid et al., 2024)	3.213	3.220	3.452	3.295
GPT-4o	3.558	3.561	3.740	3.619

Table 4: Zero-shot cross-platform transfer results of Ferret-UI One . For simplicity, we train and test the model only using data corresponding to the elementary tasks.

Training	Test - Referring					Test - Grounding				
	iPhone	iPad	AppleTV	Web	Android	iPhone	iPad	AppleTV	Web	Android
iPhone	<b>86.3</b>	68.1	31.2	45.3	71.2	<b>84.1</b>	65.2	43.1	51.7	63.1
iPad	67.5	<b>80.2</b>	40.7	51.5	63.3	64.5	<b>82.1</b>	32.1	38.5	53.8
AppleTV	29.1	45.1	<b>79.3</b>	54.2	36.4	33.7	41.2	<b>81.6</b>	52.1	29.7
Web	59.2	57.4	41.2	<b>85.5</b>	41.7	54.0	51.2	46.5	<b>87.5</b>	45.9
Android	72.5	60.7	35.7	51.2	<b>86.2</b>	66.7	48.9	29.7	44.1	<b>83.9</b>

benchmark (Chawla et al., 2024). Note, that for each data entry corresponding to the elementary and advanced tasks, it is an average across all platforms. The detailed results on each platform are provided in appendix B in Appendix. Below, we highlight a few observations.

- Ferret-UI One, powered by Llama-3-8B, delivers the best results across most metrics. It achieves the highest GPT-4o score on advanced tasks, with a notable 89.73, surpassing Ferret-UI by 43.92 points and GPT-4o by 12.0 points. Notably, Ferret-UI One with Llama-3-8B also achieves the highest IoU score on the GUIDE benchmark with 55.78, indicating superior grounding capability over the other models.
- Ferret-UI One, equipped with Vicuna-13B, also performs well, *e.g.*, achieving a strong Multi-IoU score of 41.71 on advanced tasks. Despite being six times smaller, Ferret-UI One with Gemma-2B delivers competitive performance across the board.
- In contrast, GPT-4o struggles with fine-grained UI understanding, as shown by its low referring (56.47) and grounding (12.14) scores in the elementary tasks. Its Multi-IoU and IoU scores in advanced tasks and the GUIDE benchmark are also low.

Overall, the results demonstrate the versatility of Ferret-UI One in handling UI understanding tasks across different platforms.

**Results on GUI-World.** To further demonstrate the zero-shot performance of using Ferret-UI One out of the box, we further test our model on the recently released GUI-World benchmark (Chen et al., 2024). Results are summarized in Table 3. Clearly, Ferret-UI One does not overfit to the training data, and can generalize well to the test data in the wild. Notably, Ferret-UI One outperforms GUI-Vid (Chen et al., 2024), a model developed in the GUI-World paper, on supported platforms including iOS, Android and Webpages.

### 4.3 ABLATION STUDY

**Cross-Platform Transferability.** In Table 4, we evaluate the zero-shot platform (domain) transfer capability of the Ferret-UI One model by training it on elementary tasks from one platform and testing it on other platforms. These results provide insights into how well the model generalizes across

Table 5: Ablation results of the architecture and dataset improvements of Ferret-UI One w.r.t. Ferret-UI-anyRes (You et al., 2024), *i.e.*, the high-resolution version of Ferret-UI equipped with the any-resolution module. **iPhone v1** refers to the dataset on the iPhone platform originally used by Ferret-UI, while **iPhone v2** is the data used by Ferret-UI One. Models are evaluated on the advanced tasks.

Training Data	Model	Test Data			
		iPhone v1		iPhone v2	
		GPT4o Score	Multi-IoU	GPT4o Score	Multi-IoU
iPhone v1	Ferret-UI-anyRes	91.3	36.89	68.3	27.13
	Ferret-UI One	93.7 (+2.4)	37.12 (+0.23)	70.2 (+1.9)	28.21 (+1.08)
iPhone v2	Ferret-UI-anyRes	86.2	35.89	85.97	39.81
	Ferret-UI One	88.1 (+1.9)	36.43 (+0.54)	89.7 (+3.73)	41.73 (+1.92)

different platforms. We observe similar performance patterns across both referring and grounding tasks. Specifically,

- iPhone transfers well to iPad and Android platforms on both tasks, achieving at least 68.1 and 71.2 scores on referring tasks and 65.2 and 63.1 scores on grounding tasks due to its diverse screenshot contents (because of the large screenshot number) and similar resolutions and aspect ratios with other two platforms. iPad and Android can also transfer fairly well to the iPhone domain, all achieving around 65 scores.
- AppleTV and Web do not transfer very well to the mobile domains, including iPhone, Android, and iPad, achieving the highest referring score of 59.2 and grounding score of 54.0, possibly because they are mostly landscape screenshots, which is in contrast to the mostly portrait screenshots in mobile platforms. Models trained on other domains all achieve poor performance on the AppleTV test domain, with the highest score of around 40 on both kinds of tasks, which is reasonable since there is a large domain gap in terms of AppleTV’s content distribution compared to other domains.

The results suggest that (i) iPhone, iPad and Android have similar content distribution, which help them generalize to each other; (ii) models trained on more diverse contents (*e.g.*, around 100k iPhone data) can generalize better to other platforms; (iii) platforms with similar resolution and aspect ratios may transfer better to each other; and (iv) good transferability among some of these platforms contributes to the cross-platform performance of Ferret-UI One.

**Ablation on Architecture and Dataset Improvements.** Table 5 presents a comparison between the Ferret-UI and Ferret-UI One model trained on different versions of the iPhone dataset. Models are evaluated on the test set corresponding to advanced tasks. The results indicate that both the architectural enhancements, particularly the adaptive  $N$ -gridding, and the improved dataset (iPhone v2) contribute to performance gains.

Specifically, when evaluated on the iPhone v1 test set, Ferret-UI One shows a slight improvement over Ferret-UI, with the GPT4o score increased from 91.3 to 93.7 and the Multi-IoU score increased from 36.89 to 37.12. However, the improvements are more pronounced on the iPhone v2 dataset. Here, Ferret-UI One achieves a GPT4o score of 89.7, outperforming Ferret-UI’s 85.97, along with a substantial Multi-IoU score boost from 39.81 to 41.73. These results suggest that while both architecture and dataset enhancements contribute to overall performance, the new dataset plays a more significant role in driving improvements, particularly on more challenging tasks.

## 5 CONCLUSIONS

In this paper, we presented Ferret-UI One, a multimodal large language model designed to improve UI understanding and interaction across diverse platforms. With multi-platform support, high-resolution image encoding with adaptive gridding, and improved data generation, Ferret-UI One outperforms Ferret-UI on all tested benchmarks. The model demonstrates strong zero-shot transferability across platforms, establishing Ferret-UI One as a solid foundation for universal UI understanding. Future work will focus on incorporating additional platform types and building a generalist agent for universal UI navigation.

## REFERENCES

- 540  
541  
542 Kirosos Ataallah, Xiaoqian Shen, Eslam Abdelrahman, Essam Sleiman, Deyao Zhu, Jian Ding, and  
543 Mohamed Elhoseiny. Minigt4-video: Advancing multimodal llms for video understanding with  
544 interleaved visual-textual tokens. *arXiv preprint arXiv:2404.03413*, 2024.
- 545 Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang  
546 Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities.  
547 *arXiv preprint arXiv:2308.12966*, 2023.
- 548 Rajat Chawla, Adarsh Jha, Muskaan Kumar, Mukunda NS, and Ishaan Bhola. Guide: Graphical  
549 user interface data for execution. *arXiv preprint arXiv:2404.16048*, 2024.
- 550  
551 Dongping Chen, Yue Huang, Siyuan Wu, Jingyu Tang, Liuyi Chen, Yilin Bai, Zhigang He, Chenlong  
552 Wang, Huichi Zhou, Yiqiang Li, et al. Gui-world: A dataset for gui-oriented multimodal llm-  
553 based agents. *arXiv preprint arXiv:2406.10819*, 2024.
- 554 Xingyu Chen, Zihan Zhao, Lu Chen, Danyang Zhang, Jiabao Ji, Ao Luo, Yuxuan Xiong, and  
555 Kai Yu. Websrc: A dataset for web-based structural reading comprehension. *arXiv preprint*  
556 *arXiv:2101.09465*, 2021.
- 557 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng,  
558 Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An  
559 open-source chatbot impressing gpt-4 with 90%\* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.
- 560  
561  
562 Biplab Deka, Zifeng Huang, Chad Franzen, Joshua Hibschan, Daniel Afegan, Yang Li, Jeffrey  
563 Nichols, and Ranjitha Kumar. Rico: A mobile app dataset for building data-driven design ap-  
564 plications. In *Proceedings of the 30th annual ACM symposium on user interface software and*  
565 *technology*, pp. 845–854, 2017.
- 566 Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su.  
567 Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing*  
568 *Systems*, 36, 2024.
- 569 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
570 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.  
571 *arXiv preprint arXiv:2407.21783*, 2024.
- 572 Difei Gao, Lei Ji, Zechen Bai, Mingyu Ouyang, Peiran Li, Dongxing Mao, Qinchen Wu, Weichen  
573 Zhang, Peiyi Wang, Xiangwu Guo, et al. Assistgui: Task-oriented desktop graphical user interface  
574 automation. *arXiv preprint arXiv:2312.13108*, 2023.
- 575  
576 Wenyi Hong, Weihang Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan  
577 Wang, Yuxiao Dong, Ming Ding, et al. Cogagent: A visual language model for gui agents. *arXiv*  
578 *preprint arXiv:2312.08914*, 2023.
- 579 Peng Jin, Ryuichi Takanobu, Wancai Zhang, Xiaochun Cao, and Li Yuan. Chat-univi: Unified visual  
580 representation empowers large language models with image and video understanding. In *CVPR*,  
581 2024.
- 582 Raghav Kapoor, Yash Parag Butala, Melisa Russak, Jing Yu Koh, Kiran Kamble, Waseem Alshikh,  
583 and Ruslan Salakhutdinov. Omniact: A dataset and benchmark for enabling multimodal generalist  
584 autonomous agents for desktop and web. *arXiv preprint arXiv:2402.17553*, 2024.
- 585  
586 Gang Li and Yang Li. Spotlight: Mobile ui understanding using vision-language models with a  
587 focus, 2023.
- 588 Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen,  
589 Ping Luo, et al. Mvbench: A comprehensive multi-modal video understanding benchmark. In  
590 *CVPR*, 2024.
- 591  
592 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee.  
593 Llava-next: Improved reasoning, ocr, and world knowledge, January 2024a. URL <https://llava-vl.github.io/blog/2024-01-30-llava-next/>.

- 594 Junpeng Liu, Yifan Song, Bill Yuchen Lin, Wai Lam, Graham Neubig, Yuanzhi Li, and Xiang  
595 Yue. Visualwebbench: How far have multimodal llms evolved in web page understanding and  
596 grounding? *arXiv preprint arXiv:2404.05955*, 2024b.
- 597  
598 Xiao Liu, Tianjie Zhang, Yu Gu, Iat Long Iong, Yifan Xu, Xixuan Song, Shudan Zhang, Hanyu Lai,  
599 Xinyi Liu, Hanlin Zhao, et al. Visualagentbench: Towards large multimodal models as visual  
600 foundation agents. *arXiv preprint arXiv:2408.06327*, 2024c.
- 601  
602 Kaixin Ma, Hongming Zhang, Hongwei Wang, Xiaoman Pan, Wenhao Yu, and Dong Yu. Laser:  
603 Llm agent with state-space exploration for web navigation. *arXiv preprint arXiv:2309.08172*,  
604 2023.
- 605  
606 Yichen Pan, Dehan Kong, Sida Zhou, Cheng Cui, Yifei Leng, Bing Jiang, Hangyu Liu, Yanyi Shang,  
607 Shuyan Zhou, Tongshuang Wu, et al. Webcanvas: Benchmarking web agents in online environ-  
608 ments. *arXiv preprint arXiv:2406.12373*, 2024.
- 609  
610 Christopher Rawles, Sarah Clinckemaille, Yifan Chang, Jonathan Waltz, Gabrielle Lau, Mary-  
611 beth Fair, Alice Li, William Bishop, Wei Li, Folawiyo Campbell-Ajala, et al. Androidworld: A  
612 dynamic benchmarking environment for autonomous agents. *arXiv preprint arXiv:2405.14573*,  
613 2024a.
- 614  
615 Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy Lillicrap. An-  
616 droidinthewild: A large-scale dataset for android device control. *Advances in Neural Information  
617 Processing Systems*, 36, 2024b.
- 618  
619 Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-  
620 baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gem-  
621 ini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint  
622 arXiv:2403.05530*, 2024.
- 623  
624 Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhu-  
625 patiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma  
626 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
- 627  
628 Daniel Toyama, Philippe Hamel, Anita Gergely, Gheorghe Comanici, Amelia Glaese, Zafarali  
629 Ahmed, Tyler Jackson, Shibl Mourad, and Doina Precup. Androidenv: A reinforcement learning  
630 platform for android. *arXiv preprint arXiv:2105.13231*, 2021.
- 631  
632 Junyang Wang, Haiyang Xu, Jiabo Ye, Ming Yan, Weizhou Shen, Ji Zhang, Fei Huang, and Ji-  
633 tao Sang. Mobile-agent: Autonomous multi-modal mobile device agent with visual perception,  
634 2024a.
- 635  
636 Luyuan Wang, Yongyu Deng, Yiwei Zha, Guodong Mao, Qinmin Wang, Tianchen Min, Wei Chen,  
637 and Shoufa Chen. Mobileagentbench: An efficient and user-friendly benchmark for mobile llm  
638 agents. *arXiv preprint arXiv:2406.08184*, 2024b.
- 639  
640 Jason Wu, Siyan Wang, Siman Shen, Yi-Hao Peng, Jeffrey Nichols, and Jeffrey Bigham. Webui: A  
641 dataset for enhancing visual ui understanding with web semantics. *ACM Conference on Human  
642 Factors in Computing Systems (CHI)*, 2023.
- 643  
644 Zhiyong Wu, Chengcheng Han, Zichen Ding, Zhenmin Weng, Zhoumianze Liu, Shunyu Yao, Tao  
645 Yu, and Lingpeng Kong. Os-copilot: Towards generalist computer agents with self-improvement.  
646 *arXiv preprint arXiv:2402.07456*, 2024.
- 647  
648 Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing  
649 Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, et al. Osworld: Benchmarking multimodal  
650 agents for open-ended tasks in real computer environments. *arXiv preprint arXiv:2404.07972*,  
651 2024.
- 652  
653 Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chunyuan Li, and Jianfeng Gao. Set-of-mark  
654 prompting unleashes extraordinary visual grounding in gpt-4v. *arXiv preprint arXiv:2310.11441*,  
655 2023a.

648 Zhao Yang, Jiaxuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. Appagent:  
649 Multimodal agents as smartphone users. *arXiv preprint arXiv:2312.13771*, 2023b.  
650

651 Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable  
652 real-world web interaction with grounded language agents. *Advances in Neural Information Pro-*  
653 *cessing Systems*, 35:20744–20757, 2022.

654 Haoxuan You, Haotian Zhang, Zhe Gan, Xianzhi Du, Bowen Zhang, Zirui Wang, Liangliang Cao,  
655 Shih-Fu Chang, and Yinfei Yang. Ferret: Refer and ground anything anywhere at any granularity,  
656 2023.  
657

658 Keen You, Haotian Zhang, Eldon Schoop, Floris Weers, Amanda Swearngin, Jeffrey Nichols, Yinfei  
659 Yang, and Zhe Gan. Ferret-ui: Grounded mobile ui understanding with multimodal llms. *arXiv*  
660 *preprint arXiv:2404.05719*, 2024.

661 Chaoyun Zhang, Liqun Li, Shilin He, Xu Zhang, Bo Qiao, Si Qin, Minghua Ma, Yu Kang, Qingwei  
662 Lin, Saravan Rajmohan, et al. Ufo: A ui-focused agent for windows os interaction. *arXiv preprint*  
663 *arXiv:2402.07939*, 2024a.

664 Danyang Zhang, Hongshen Xu, Zihan Zhao, Lu Chen, Ruisheng Cao, and Kai Yu. Mobile-env:  
665 an evaluation platform and benchmark for llm-gui interaction. *arXiv preprint arXiv:2305.08144*,  
666 2023.  
667

668 Haotian Zhang, Haoxuan You, Philipp Dufter, Bowen Zhang, Chen Chen, Hong-You Chen, Tsu-Jui  
669 Fu, William Yang Wang, Shih-Fu Chang, Zhe Gan, et al. Ferret-v2: An improved baseline for  
670 referring and grounding with large language models. *arXiv preprint arXiv:2404.07973*, 2024b.

671 Longtao Zheng, Zhiyuan Huang, Zhenghai Xue, Xinrun Wang, Bo An, and Shuicheng Yan.  
672 Agentstudio: A toolkit for building general virtual agents, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2403.17918)  
673 [abs/2403.17918](https://arxiv.org/abs/2403.17918).  
674

675 Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng,  
676 Tianyue Ou, Yonatan Bisk, Daniel Fried, et al. Webarena: A realistic web environment for build-  
677 ing autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023.  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

Table 6: The performance breakdown of Ferret-UI One with Llama-3-8B backbone. Note that we only report GPT-4o evaluation scores and omit the multi-IoU scores for advanced tasks for simplicity.

Task type	Task	Test Domain				
		iPhone	iPad	AppleTV	Web	Android
Refer	OCR	75.3	69.3	74.3	80.5	74.5
	Widget Classify	79.1	78.7	82.5	83.6	82.0
	Tapperbility	89.2	86.0	-	84.3	85.6
Ground	Widget Listing	76.7	74.9	71.6	79.7	76.8
	Find Text	85.2	84.0	82.2	90.2	82.2
	Find Widget	88.6	86.8	87.7	87.8	86.4
Advanced Tasks	Comprehensive	89.0	86.5	84.1	83.7	85.6
	Perception	94.1	87.5	86.6	93.2	92.1
	Interaction	93.5	92.7	86.1	96.4	94.7

## A DETAILED GPT-4O REQUIREMENTS FOR GENERATING ADVANCED TASKS

In this section, we elaborate how we prompt GPT-4o with how to generate advanced tasks. In particular, we have following requirements when generating each type of advanced tasks:

- Comprehensive Description:** Provide a one-sentence description of the overall functionality of the UI page shown in the screenshot. Then, describe the screenshot in detail by dividing it into several areas/groups and explaining the functionality of each area/group.
- Multi-Round Perception QA:**
  - Basic Perceptual Understanding:** Interpret the content of referred widgets, ground the positions of the widgets based on the widget descriptions, and describe the functionality and status (enabled, disabled, selected, hovered) of each widget.
  - Contextual Awareness:** Understand the context in which UI widgets are presented, including the relationships between different UI components and how they contribute to the overall user experience.
  - Layout and Hierarchy Recognition:** Recognize how widgets are grouped and nested within the UI layout.
- Multi-Round Interaction QA:**
  - Generate instructions for user-centered interaction. For example, “please help me confirm submission” instead of “please click on [Box0] button”.
  - Identify and interact with buttons, links, icons, scrollbars, and toggles, and recognize widgets that cannot be interacted with, such as static text and background images.
  - Identify and interact with input fields, dropdowns, checkboxes, and radio buttons, as well as their current states (e.g., checked, unchecked, filled, empty) in the context of specific interaction tasks.
  - Predict current step for potential navigation paths based on currently visible widgets like tabs, menus and breadcrumbs controls.

## B DETAILED RESULTS ON ELEMENTARY AND ADVANCED TASKS

Table 6 showcases the performance on 9 subtasks  $\times$  5 platforms of the Ferret-UI One model. We empirically find that OCR, widget listing and comprehensive description are relatively difficult tasks compared to other tasks within the same task type.

756 C ADVANCED TASK EXAMPLES  
757  
758  
759  
760  
761

762 Comprehensive Description

763  
764 **OVERALL:** This UI page is for adjusting settings related to downloads, language, news,  
765 and notifications.

766 **DETAIL:**

- 767 • **Top Navigation Bar:** Contains menu options such as Home, Account, Settings (currently  
768 selected), and Downloads, allowing users to navigate between different sections.
- 769 • **Games Section:** Features a checkbox labeled "Download over WiFi only," which is cur-  
770 rently checked, indicating that downloads are restricted to WiFi connections.
- 771 • **Language Section:** Offers a dropdown menu to select the language, currently set to En-  
772 glish, with a pink arrow indicating more options.
- 773 • **Title Screen Section:** Includes a checkbox for "Show Sago News," which is checked,  
774 meaning news updates will be displayed on the title screen.
- 775 • **Notifications Section:** Contains a checkbox for "Display notifications," which is  
776 unchecked, allowing the user to enable or disable notifications.
- 777 • **Close Button:** Located in the top right corner, represented by an 'X', used to close the  
778 settings page.  
779

780  
781  
782  
783  
784  
785  
786 Multi-Round Perception QA

787 **Q:** *What does the "SERVEWARE" text represent?*

788 **A:** The text "SERVEWARE" [718, 72, 948, 108] is likely a category label indicating that the  
789 items displayed are related to serving ware.  
790

791 **Q:** *Can you identify any product names and their corresponding prices?*

792 **A:** Yes, the product "ARV BRÖLLOP" [35, 847, 233, 907] is priced at \$12.99 [10, 6, 112,  
793 32] and "IKEA 365+" [609, 787, 809, 823] is priced at \$2.99 [1174, 1336, 1627, 2119].  
794

795 **Q:** *What is the function of the shopping cart icon?*

796 **A:** The shopping cart icon [630, 2157, 690, 2217] is typically used to view the items in your  
797 cart or proceed to checkout.

798 **Q:** *How is the product rating displayed?*

799 **A:** The product rating is displayed with star symbols and a numerical rating. For example,  
800 "ARV BRÖLLOP" has a rating of 4.9 next to the stars [[38, 956, 100, 986], [248, 1004, 294,  
801 1030]].  
802

803 **Q:** *What does the "Filter" button do?*

804 **A:** The "Filter" button [1564, 170, 1630, 198] is likely used to apply specific criteria to  
805 narrow down the product search results.

806 **Q:** *Are there any navigation elements on the screen?*

807 **A:** Yes, there seems to be a navigation bar at the bottom with icons for different sections,  
808 like a house icon for home [611, 268, 1057, 714] and a pencil icon [1320, 2157, 1380, 2217]  
809 which might be for editing or creating a list.

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

**Multi-Round Interaction QA**

**Q:** *How can I refresh the content shown?*  
**A:** You can tap on the refresh icon [272, 736, 424, 882] to reload the content.

**Q:** *What does the "General" tag refer to?*  
**A:** The "General" tag [1074, 900, 1443, 926] categorizes the content related to general events, like the Caltech Beaver Name Reveal.

**Q:** *How can I view more recent activities?*  
**A:** You can tap the "Recent" section label [90, 598, 202, 630] to view more recent activities.

**Q:** *What is the logo displayed on this page?*  
**A:** The logo at the top of the page is the SCIAN Network logo [198, 56, 444, 128] representing the network conducting the event.

**Q:** *How do I check the institution hosting the event?*  
**A:** The event is hosted by the California Institute of Technology, as shown in the text [11, 499, 483, 542] below the event details.

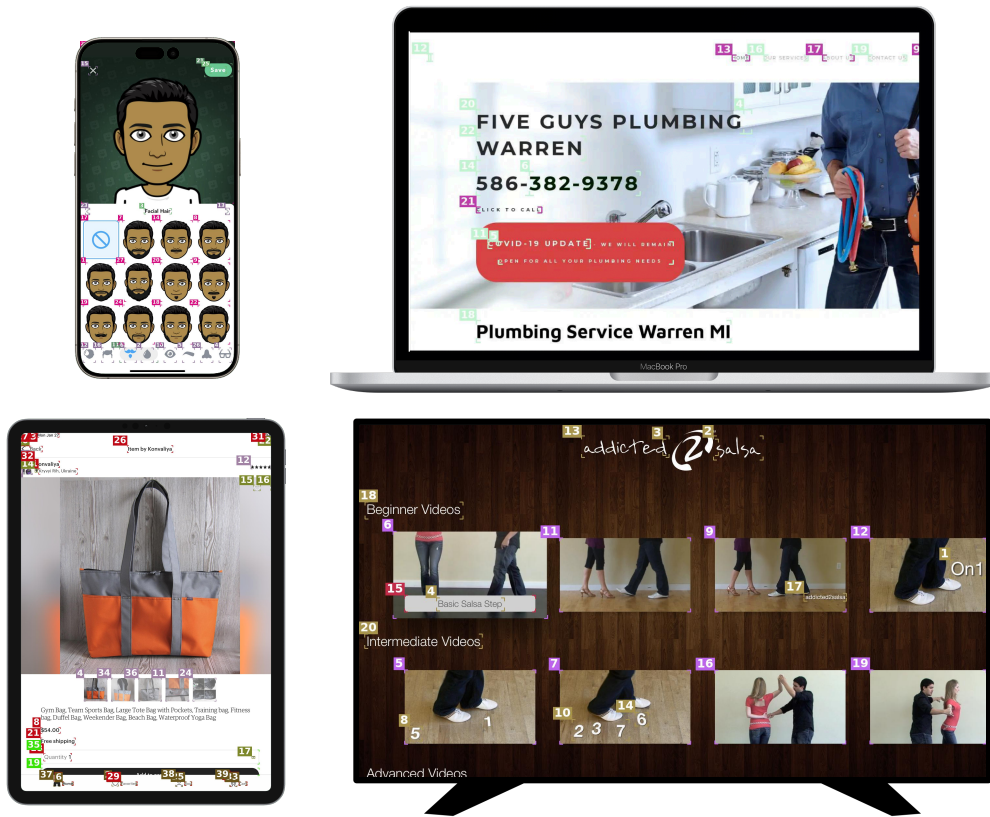


Figure 5: Examples of visual prompting using GPT-4o to generate task data for Multi-Round Perception QA and Multi-Round Interaction QA. Each UI widget is annotated with a corner-style bounding box, where only the corners of the widget are highlighted by small lines, leaving the rest of the box open. This minimalistic bounding style is accompanied by a unique number tag placed near one of the corners, making it easy to identify and reference specific UI widgets for further interaction or perception analysis.



## D RESOLUTION STATISTICS

In this section, we present the resolution statistics for images collected from various platforms, categorized by device type and resolution.

**iPhone** For iPhone, the following resolutions were observed:

- 828x1792: 83,250 images
- 1125x2436: 6,055 images
- 1792x828: 4,686 images
- 2436x1125: 104 images

**iPad** For iPad, the observed resolutions are:

- 2224x1668: 4,829 images
- 1668x2224: 14,312 images
- 1242x2208: 19 images

**AppleTV** For AppleTV, the resolution of 1920x1080 accounted for 16,152 images.

**WebUI** For the webpage data from WebUI dataset, each screenshot contains the following 6 versions of resolutions, we randomly pick one resolution for each screenshot, resulting in following resolution distribution:

- 1280x720: 53,500 images
- 1366x768: 53,500 images
- 1536x864: 53,500 images
- 1920x1080: 53,500 images
- 2048x2732: 53,500 images
- 1170x2532: 53,500 images

**Android** For Android data from RICO dataset, the following resolutions were observed:

- 540x960: 14,092 images
- 1080x1920: 52,102 images
- 1920x1080: 55 images
- 960x540: 12 images

## E LABEL STATISTICS AND MAPPING RESULTS OF ORIGINAL DATA ACROSS PLATFORMS

In this section, we demonstrate the label statistics of original data across platforms and their mapping results into a uniform label space for better joint training. The mapping results are obtained by the combination of GPT-4 suggestions and human review. Note that “Other” label after mapping indicates the widget info will be deprecated.

### E.1 IOS LABELS (IPHONE + IPAD)

- **Text:** 867450, Text
- **Icon:** 331390, Icon
- **Container:** 236801, Container
- **Tab:** 203839, TabBarItem

- 918 • **Picture**: 187222, Picture
- 919 • **TextField**: 59082, TextField
- 920
- 921 • **SegmentedControl**: 60310, SegmentedControl
- 922 • **Other**: 17084, Other
- 923 • **Checkbox**: 12203, Checkbox
- 924 • **PageControl**: 9074, PageControl
- 925 • **Toggle**: 8327, Toggle
- 926 • **Slider**: 4396, Slider
- 927 • **SkipScreen**: 2352, Other
- 928
- 929

## 930 E.2 APPLETV LABELS

- 931
- 932 • **Checkbox (Not Selected)**: 1, Checkbox
- 933 • **Checkbox (Selected)**: 1, Checkbox
- 934 • **Container**: 1, Container
- 935 • **Dialog**: 1, Dialog
- 936 • **Icon**: 1, Icon
- 937 • **PageControl**: 1, PageControl
- 938 • **Picture**: 1, Picture
- 939 • **SegmentedControl**: 1, SegmentedControl
- 940 • **Slider**: 1, Slider
- 941 • **Tab**: 1, TabBar
- 942 • **Text**: 1, Text
- 943 • **TextField**: 1, TextField
- 944 • **Toggle (Not Selected)**: 1, Toggle
- 945 • **Toggle (Selected)**: 1, Toggle
- 946 • **Other**: 1, OtherUI
- 947 • **SkipScreen**: 1, Other
- 948
- 949
- 950
- 951

## 952 E.3 ANDROID LABELS

- 953
- 954 • **Toolbar**: 34854, Other
- 955 • **Icon**: 179016, Icon
- 956 • **Text**: 455075, Text
- 957 • **List Item**: 152165, Container
- 958 • **Card**: 16901, Container
- 959 • **Text Button**: 138834, Button
- 960 • **Image**: 217511, Picture
- 961 • **Radio Button**: 5420, Checkbox
- 962 • **Input**: 21408, TextField
- 963 • **Advertisement**: 13338, Other
- 964 • **Web View**: 30788, Other
- 965 • **Pager Indicator**: 4144, PageControl
- 966 • **Background Image**: 6785, Picture
- 967 • **Slider**: 2016, Slider
- 968 • **Drawer**: 6642, Other
- 969
- 970
- 971

- 972 • **Multi-Tab**: 4185, Tab
- 973 • **Modal**: 3959, Dialog
- 974 • **Button Bar**: 721, Other
- 975 • **Checkbox**: 4243, Checkbox
- 976 • **Bottom Navigation**: 524, Other
- 977 • **Number Stepper**: 427, Other
- 978 • **Date Picker**: 291, Other
- 979 • **On/Off Switch**: 2103, Toggle
- 980 • **Map View**: 1512, Other
- 981 • **Video**: 562, Other
- 982
- 983
- 984
- 985 E.4 WEBPAGE LABELS
- 986
- 987 • **StaticText**: 10397304, Text
- 988 • **link**: 6613952, Button
- 989 • **listitem**: 4142046, Other
- 990 • **generic**: 3981172, Other
- 991 • **img**: 1222637, Picture
- 992 • **heading**: 1109202, Other
- 993 • **paragraph**: 1028955, Other
- 994 • **list**: 760601, Other
- 995 • **LineBreak**: 377966, Other
- 996 • **Section**: 216403, Other
- 997 • **button**: 197497, Button
- 998 • **gridcell**: 178174, Other
- 999 • **ListMarker**: 152519, Other
- 1000 • **LayoutTableCell**: 144391, Other
- 1001 • **banner**: 118724, Other
- 1002 • **navigation**: 106019, Other
- 1003 • **textbox**: 102802, TextField
- 1004 • **emphasis**: 97256, Other
- 1005 • **DescriptionListDetail**: 96085, Other
- 1006 • **DescriptionListTerm**: 88250, Other
- 1007 • **DescriptionList**: 83379, Other
- 1008 • **separator**: 79504, Other
- 1009 • **contentinfo**: 73695, Other
- 1010 • **LayoutTableRow**: 71033, Other
- 1011 • **time**: 66272, Other
- 1012 • **row**: 65959, Other
- 1013 • **article**: 57367, Other
- 1014 • **LabelText**: 54504, Other
- 1015 • **LayoutTable**: 51826, Other
- 1016 • **HeaderAsNonLandmark**: 45344, Other
- 1017 • **strong**: 42913, Other
- 1018 • **insertion**: 40704, Other
- 1019
- 1020
- 1021
- 1022
- 1023
- 1024
- 1025

- 1026 • **superscript**: 37056, Other
- 1027 • **RootWebArea**: 31708, Other
- 1028 • **complementary**: 27495, Other
- 1029 • **Iframe**: 24465, Other
- 1030 • **FooterAsNonLandmark**: 19971, Other
- 1031 • **combobox**: 19960, Toggle
- 1032 • **main**: 17984, Other
- 1033 • **figure**: 17704, Other
- 1034 • **search**: 16930, Other
- 1035 • **columnheader**: 15141, Other
- 1036 • **menuitem**: 14390, Other
- 1037 • **searchbox**: 13544, TextField
- 1038 • **group**: 10545, Other
- 1039 • **table**: 10281, Other
- 1040 • **checkbox**: 10148, Checkbox
- 1041 • **option**: 9397, Other
- 1042 • **Pre**: 7503, Other
- 1043 • **graphics-symbol**: 6991, Other
- 1044 • **radio**: 5385, Checkbox
- 1045 • **rowheader**: 4685, Other
- 1046 • **region**: 4488, Other
- 1047 • **blockquote**: 3989, Other
- 1048 • **dialog**: 3673, Dialog
- 1049 • **Canvas**: 3664, Other
- 1050 • **tab**: 3316, Tab
- 1051 • **Abbr**: 3141, Other
- 1052 • **SvgRoot**: 2969, Other
- 1053 • **PluginObject**: 2918, Other
- 1054 • **Figcaption**: 2847, Other
- 1055 • **IframePresentational**: 2590, Other
- 1056 • **EmbeddedObject**: 2516, Other
- 1057 • **menu**: 2172, Other
- 1058 • **code**: 1900, Other
- 1059 • **Video**: 1682, Other
- 1060 • **Legend**: 1651, Other
- 1061 • **Details**: 1513, Other
- 1062 • **alert**: 1466, Other
- 1063 • **menubar**: 1452, Other
- 1064 • **status**: 1363, Other
- 1065 • **DisclosureTriangle**: 1338, Other
- 1066 • **tabpanel**: 1160, Other
- 1067 • **tablist**: 1122, Other
- 1068 • **listbox**: 1064, Other
- 1069 • **form**: 1050, Other

- 1080 • **alertydialog**: 1007, Dialog
- 1081 • **mark**: 1003, Other
- 1082 • **slider**: 941, Slider
- 1083 • **progressbar**: 903, Other
- 1084 • **treeitem**: 845, Other
- 1085 • **spinbutton**: 726, Other
- 1086 • **document**: 625, Other
- 1087 • **subscript**: 519, Text
- 1088 • **note**: 483, Other
- 1089 • **deletion**: 429, Other
- 1090 • **application**: 273, Other
- 1091 • **caption**: 262, Other
- 1092 • **rowgroup**: 247, Other
- 1093 • **switch**: 233, Toggle
- 1094 • **grid**: 216, Other
- 1095 • **log**: 181, Other
- 1096 • **toolbar**: 169, Other
- 1097 • **math**: 167, Other
- 1098 • **tooltip**: 151, Other
- 1099 • **radiogroup**: 151, Other
- 1100 • **Audio**: 132, Other
- 1101 • **meter**: 122, Other
- 1102 • **Ruby**: 88, Other
- 1103 • **doc-noteref**: 87, Other
- 1104 • **timer**: 86, Other
- 1105 • **menuitemradio**: 84, Other
- 1106 • **tree**: 71, Other
- 1107 • **definition**: 33, Other
- 1108 • **Date**: 30, Other
- 1109 • **graphics-object**: 29, Other
- 1110 • **doc-toc**: 23, Other
- 1111 • **doc-backlink**: 22, Other
- 1112 • **ColorWell**: 14, Other
- 1113 • **doc-subtitle**: 13, Other
- 1114 • **marquee**: 12, Other
- 1115 • **directory**: 9, Other
- 1116 • **menuitemcheckbox**: 8, Other
- 1117 • **doc-endnote**: 4, Other
- 1118 • **doc-endnotes**: 3, Other
- 1119 • **feed**: 3, Other
- 1120 • **doc-abstract**: 2, Other
- 1121 • **treegrid**: 1, Other
- 1122 • **DateTime**: 1, Other
- 1123
- 1124
- 1125
- 1126
- 1127
- 1128
- 1129
- 1130
- 1131
- 1132
- 1133