



# Ferret-UI: Grounded Mobile UI Understanding with Multimodal LLMs

Keen You, Haotian Zhang, Eldon Schoop, Floris Weers,  
Amanda Swearngin, Jeffrey Nichols, Yinfei Yang, and Zhe Gan

Apple

{k\_you,haotian\_zhang2,eldon,fweers,aswearngin,jwnichols,yinfeiy,zhe.gan}@apple.com

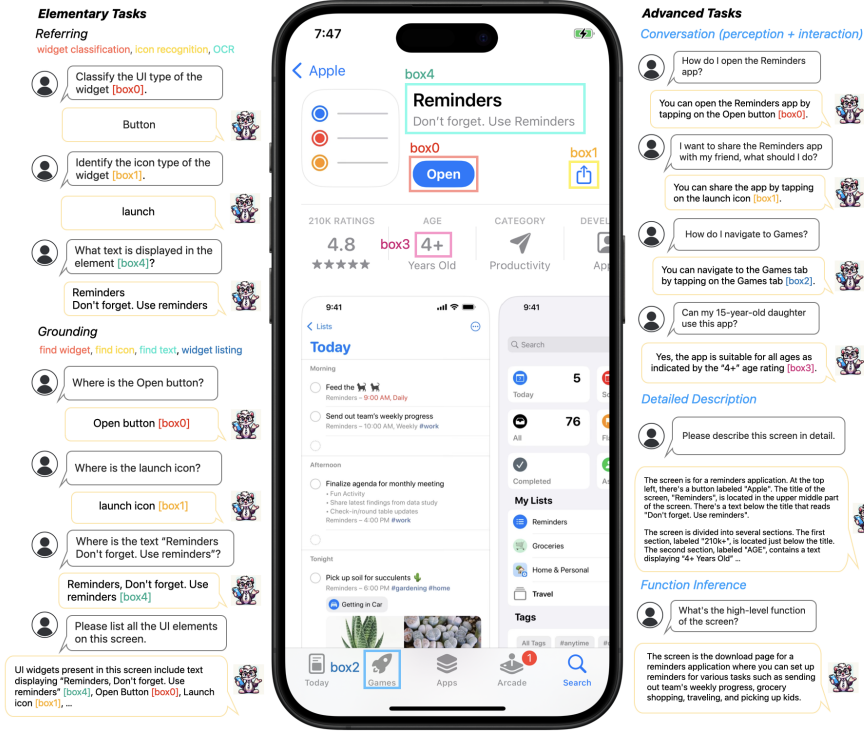
**Abstract.** The recent advancements in multimodal large language models (MLLMs) have been noteworthy, yet, these general-domain MLLMs often fall short in their ability to comprehend and interact effectively with user interface (UI) screens. In this paper, we construct Ferret-UI, a new MLLM tailored for enhanced understanding of mobile UI screens, equipped with *referring*, *grounding*, and *reasoning* capabilities. We meticulously gathered training samples from an extensive range of fundamental UI tasks, such as *icon recognition*, *find text*, and *widget listing*. These samples are formatted for instruction-following with region annotations to facilitate precise referring and grounding. Moreover, to augment the model’s reasoning ability, we compile a dataset for advanced tasks inspired by Ferret, but with a focus on mobile screens. This methodology enables the training of Ferret-UI, a model that exhibits outstanding comprehension of UI screens and the ability to execute open-ended instructions, thereby facilitating UI operations. To rigorously evaluate its capabilities, we establish a comprehensive benchmark encompassing the aforementioned tasks. Ferret-UI not only outstrips most open-source UI MLLMs in performance but also achieves parity with GPT-4V, marking a significant advancement in the field.

**Keywords:** UI Understanding · Multimodal Large Language Model (MLLM)

## 1 Introduction

Using mobile applications has become an important part of daily life. Particularly, people use different mobile applications to achieve personal goals such as information search, reservation, entertainment, and directions. In this usage, we inspect the current screen visually, and perform the desired actions based on our goals. Automating this process of perception and interaction has the potential to help users achieve their goals with relative ease. Moreover, it is also a valuable building block for accessibility [13], multi-step UI navigation [18,40,48], app testing [2,29], usability studies [22], and many others.

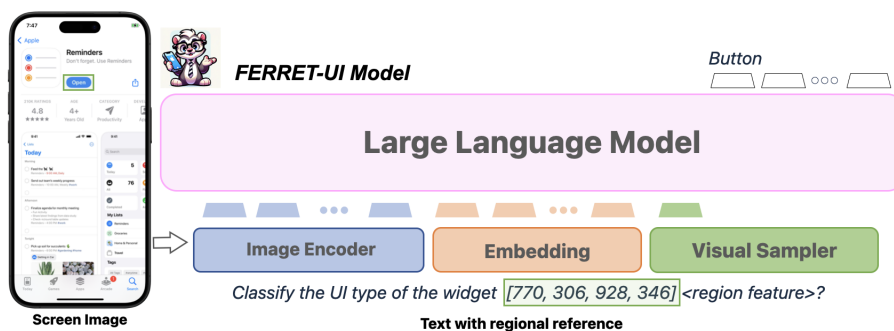
To facilitate the seamless automation of perception and interaction within user interfaces, a sophisticated system endowed with a set of key capabilities



**Fig. 1:** Ferret-UI is able to perform *referring* tasks (e.g., *widget classification*, *icon recognition*, *OCR*) and *grounding* tasks (e.g., *find widget*, *find icon*, *find text*, *widget listing*) on mobile UI screens. These elementary tasks enrich the model with both visual and spatial knowledge, enabling it to distinguish among different UI types at both broad and fine levels, such as between various icons or text elements. This foundational knowledge is crucial for performing more advanced tasks. Specifically, Ferret-UI is able to not only discuss visual elements in *detailed description* and *perception conversation*, but also propose goal-oriented actions in *interaction conversation* and deduce the overall function of the screen via *function inference*.

is essential. Such a system must possess the ability to not only apprehend the entirety of a screen but also to concentrate on specific UI elements within that screen. With visual understanding as the foundation, it should further be able to map natural language instructions to corresponding actions within a given UI, execute advanced reasoning, and provide exhaustive details concerning the screens it interacts with. These requirements necessitate the development of a vision-language model adept at both referring and grounding in relation to UI screens. Here, *referring* requires the system to utilize particular regional image information in the screen input, while *grounding* involves the model’s capacity to identify and denote precise locations on the screen in its outputs.

To achieve our objectives in enhancing mobile UI understanding, we focus on three pivotal dimensions: network architecture, data creation, and establishment



**Fig. 2:** Overview of Ferret-UI architecture. A pre-trained image encoder and projection layer produces image features for the entire screen. For text with regional references, a visual sampler generates a corresponding regional continuous feature. The LLM uses the full-image representation, regional feature, and text embeddings to generate a response.

of benchmarks. For network architecture, we base our approach on Ferret [46], a Multimodal Large Language Model (MLLM) known for its strong performance to perform both referring and grounding with natural images. We posit that Ferret’s capabilities provide a solid foundation in interactive UI-centric tasks.

Using MLLMs for UI tasks is not a new concept. However, existing work have primarily focused on processing entire screens as singular inputs (*e.g.*, Pix2Struct [25], ILuvUI [21], CogAgent [18]), only supports referring tasks with one bounding box in the input (*e.g.*, Spotlight [27]), leveraging GPT-4V [44] to navigate UI screens, as seen in MM-Navigator [42], AppAgent [48], MobileAgent [40]). In contrast, Ferret-UI emerges as the inaugural MLLM designed to execute precise referring and grounding tasks specific to UI screens, while adeptly interpreting and acting upon open-ended language directives.

To train Ferret-UI, we generate data at different granularities, covering basic semantic and spatial tasks for UI primitives to advanced reasoning tasks. We first generate training samples for basic UI tasks using a template-based approach. This encompasses *referring* tasks such as *widget classification*, *icon recognition*, *OCR*, and *grounding* tasks like *find widget*, *find icon*, *find text*, and *widget listing*. These tasks are instrumental in teaching the model to understand the semantics and spatial positioning of UI elements, enabling the model to make distinctions at both a broad level (among various UI types) and a more detailed level (within specific UI types, such as icons or text). For advanced tasks, we use GPT-4 [33] to generate data, including *detailed description*, *conversation perception*, *conversation interaction*, and *function inference*. These advanced tasks prepare the model to engage in more nuanced discussions about visual components, formulate action plans with specific goals in mind, and interpret the general purpose of a screen. Fig. 1 illustrates examples of Ferret-UI’s proficiency in handling the 11 tasks ranging from basic to advanced.

To assess these capabilities, we develop a comprehensive test benchmark featuring 14 diverse mobile UI tasks in terms of referring and grounding. This

includes 3 tasks from Spotlight [27] (*screen2words*, *widget captions*, and *taperception*), and dual versions of the 11 UI tasks previously described, tailored for both iPhone and Android screens. In our evaluation process, we conduct comprehensive benchmarking of a variety of UI understanding models. This includes assessments of open-source models such as CogAgent [18] and Fuyu [5], alongside evaluations of commercial models, including GPT-4V. We observe that Ferret-UI significantly surpasses the base Ferret model, illustrating a significant domain-specific gap. Ferret-UI is also on par with GPT-4V and notably, in the context of advanced tasks, Ferret-UI surpasses the performance of both Fuyu and CogAgent.

Our contributions are summarized as follows. (i) We have conceptualized and designed a hierarchy that encompasses both elementary and advanced UI tasks, for which we have meticulously gathered training samples specifically tailored to facilitate model training. (ii) We have developed Ferret-UI, the first UI-centric MLLM that is capable of effectively executing referring, grounding, and reasoning tasks. (iii) We have constructed a comprehensive test suite of all the tasks encompassing all tasks under investigation, and through rigorous experiments and analysis, we offer insights into the model’s capabilities and limitations.

## 2 Preliminaries

**Existing MLLMs for UI.** Earlier works [6, 17, 28, 30, 37] in the area focus on studying simplified web and mobile screens. With recent advances in both LLMs and MLLMs [3, 12, 15, 19, 20, 33, 39], the approaches to many problems have been transformed, including the UI field. Many works have demonstrate the capability of using other existing MLLMs [9, 26, 31, 38, 45, 55] for UI tasks. ILuvUI [21] and Spotlight [27] concentrating on single-screen UI tasks with exploring various UI tasks by fine-tuning on GPT-generated data, and delving into UI tasks such as screen summarization and widget interaction. MobileAgent [40] and AppAgent [48] represent a different approach, utilizing MLLMs as agents for UI screen navigation, with MobileAgent employing external detection modules for action generation and AppAgent leveraging overlaid UI element IDs and screen XML files for predefined actions. CogAgent [18], built upon CogVLM [41], shifts the focus towards using only screen images for complex UI navigation, eliminating the need for UI-specific modules. These are some examples among other works that utilize LLMs [11, 16, 23, 54] and MLLMs [8, 14, 36, 42, 47, 52] in the space.

**Ferret: An interactive MLLM.** In recent investigations, there has been a growing focus on the convergence of models [7, 24, 34, 46, 49, 51] and the tasks related to dense visual perception, where these models focus on more fine-grained interactions. Ferret-UI is based on Ferret [46]. Ferret distinguishes itself from other MLLMs by excelling in spatial referring and grounding within natural images of diverse shapes and levels of detail. It stands out for its ability to interpret and interact with regions or objects, whether they are specified as points, boxes, or any free-form shapes. The model incorporates a unique hybrid representation technique that transforms specified regions into a format suitable for



Platform	Resolution	Train	Test	Task	iPhone	Android
Android	2560×1440	26,527	3,080	screen2words	-	78k
	1792×828	74,953	8,297	widget captions	-	109k
iPhone	828×1792	4,225	461	taperception	-	14k
	2436×1125	5,420	635	elementary tasks	40k×7	40k×7
	1125×2436	87	17	advanced tasks	10k×4	10k×4

(a) Number of screens by resolution. (b) Number of samples per training task.

**Table 1:** Mobile UI screen and training data statistics.

processing by the LLMs. Ferret utilized a sophisticated architecture comprising a pre-trained visual encoder (CLIP-ViT-L/14) [35] and a decoder-only language model (Vicuna [53]). At its core, a spatial-aware visual sampler is designed to adeptly manage continuous features of region shapes in different sparsity levels. We illustrate the overall architecture of Ferret/Ferret-UI in Fig. 2. These meticulous design choices enable Ferret to interactively process user inputs and deliver outputs with unparalleled accuracy and contextual relevance, highlighting its practical application in specialized UI tasks.

Ferret-UI includes a broader range of referring tasks (e.g., *OCR*, *icon recognition*, *widget classification*) and grounding tasks (e.g., *find text/icon/widget*, *widget listing*), enhancing the model’s reasoning capabilities for advanced UI interactions. Unlike previous MLLMs that require external detection modules or screen view files, our Ferret-UI model is self-sufficient, leveraging embedded knowledge acquired during training. This approach not only facilitates advanced single-screen interactions but also paves the way for new applications, such as improving accessibility. Our work aims to specialize in UI tasks, demonstrating the potential of multimodal large language models (MLLMs) to revolutionize UI understanding and interaction.

### 3 Dataset and Task Formulation

In this section, we detail the process of generating datasets for model training and evaluation. Specifically, we describe the UI detection data collection process in Section 3.1, and we outline how we create task-specific data from raw detections in Section 3.2.

#### 3.1 UI Data Collection

**UI Screens.** To build a model capable of perceiving and interacting with mobile screens, it is crucial to gather a varied collection of such screens. This study examines screens from both iPhone and Android devices.

For Android screens, we use a subset of the RICO dataset [10]. Specifically, we consider the tasks in Spotlight [27], whose data is publicly available, including



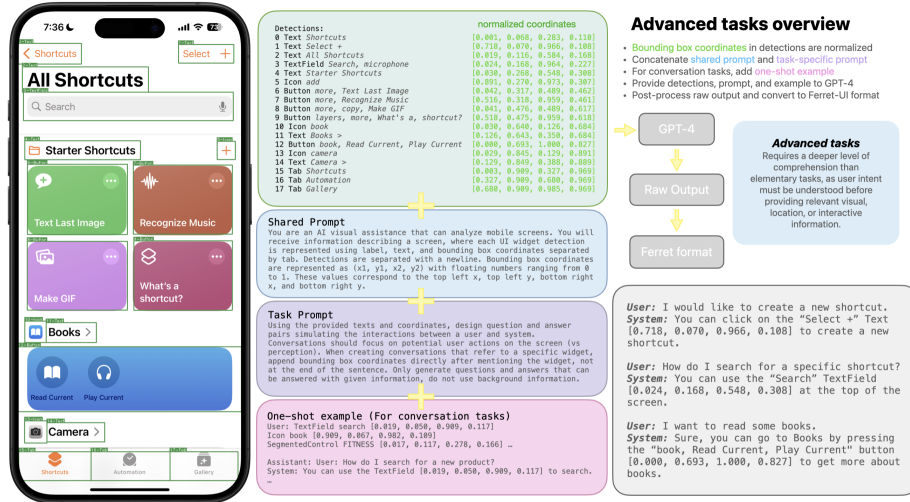
**Fig. 3: Elementary task data generation overview.** A UI detection model outputs all detected elements, with each element’s *type*, *text*, and *bounding boxes*. These detections are used to create training samples for elementary tasks. For *grounding tasks*, we use all element detections to create one sample for widget listing whereas the remaining tasks focus on one element at a time. We separate the elements into *icons*, *text*, and *non-icon/text widgets*. For each type, we create one referring sample and one grounding sample.

*screen2words*, *widgetcaptions*, and *taperception*. We aggregate unique images for each split (train and test) among the tasks to form our own data splits. In total, there are 26,527 train images and 3,080 test images.

For iPhone screens, we use the AMP dataset [50], which spans a broad spectrum of applications. A subset is randomly selected and divided into training and test splits. The iPhone screens come in various sizes, resulting in a total of 84,685 training images and 9,410 test images. The breakdown of image sizes is summarized in Tab. 1a.

**UI Screen Elements Annotation.** After collecting Android and iPhone screens, we further collect fine-grained element annotation from screens using a pre-trained pixel-based UI detection model [50]. For each detected UI element, the output includes a UI type (Button, Text, Icon, Picture, *etc.*), the corresponding bounding box, and the text displayed on it, if any, identified by the Apple Vision Framework<sup>1</sup>. We further use heuristics from Screen Recognition [50] to group individual detections into larger units, *e.g.*, multiple lines of text are merged into one group, an image is grouped with its caption, *etc.*

<sup>1</sup> <https://developer.apple.com/documentation/vision>



**Fig. 4: Advanced task data generation overview.** We first normalize bounding box coordinates from the detection outputs, then we send the detections, prompts, and optional one-shot example to GPT-4. For detailed description and function inference, we pair the generated response with a pre-selection of prompts to train Ferret-UI. For conversation tasks, we directly transform GPT-4 output to multi-turn conversations.

### 3.2 Task Formulation

This section describes how we convert the UI screens along with the associated detection data to a format that can be used to train an MLLM. We elaborate on three distinct approaches for the construction of the dataset.

**Reformatting Spotlight.** We first take *screen2words*, *widgetcaptions*, and *taperception* from the existing Spotlight tasks [27], and format them into conversational QA pairs. Specifically, GPT-3.5 Turbo is used to create a varied set of prompts from base prompts we author for respective tasks:

- **Screen2words:** *Provide a summary of this screenshot,*
- **Widget Captions:** *For the interactive element [bbox], provide a phrase that best describes its functionality,*
- **Taperception:** *Predict whether the UI element [bbox] is tappable.*

For each training example, we sample a prompt for the corresponding task, and pair it with the original source image and ground truth answer.

**Elementary Tasks.** In addition to the Spotlight tasks, we use paired screens and UI elements mentioned in Section 3.1 to generate data for novel UI tasks that rely on grounding and referring capabilities. We introduce 7 tasks using this approach, one set for each of Android and iPhone screens: *OCR*, *icon recognition*, and *widget classification* for *referring*; and *widget listing*, *find text*, *find icon*, and *find widget* for *grounding*. We define *referring tasks* as the ones with bounding

boxes in the inputs, while *grounding tasks* are the ones with bounding boxes in the outputs.

For each task, we also use GPT-3.5 Turbo to expand a base prompt to introduce variants of the task question. Details for data generation are illustrated in Fig. 3. The number of training samples for each task is summarized in Tab. 1b. The number of test samples for all tasks are 5k. In experiments, we sample from this pool of training data with different ratios to construct our training data mixture.

**Advanced Tasks.** To incorporate reasoning abilities into our model, we follow LLaVA [31], and additionally collect data of 4 more formats using GPT-4. We focus on iPhone screens for this part of the data collection, filtering our examples to those with more than 2 but fewer than 15 detections. These examples are sent together with prompts to GPT-4 to create data of the desired format—the actual images are not used. Fig. 4 illustrates the training data generation process for advanced tasks.

The four tasks are *detailed description*, *conversation perception*, *conversation interaction*, and *function inference*. Among these, we expand base prompts for detailed description and function inference to pair them with the GPT-4 response as the input data in our model training. For conversations, we provide an in-context example for GPT-4 to better follow bounding box formats in its output. From the raw GPT-4 output, we parse the bounding boxes and transform them into the correct multi-turn conversation format for our model. In total, we have created 40K valid conversations from GPT-4 generated data. More details about our data collection pipeline, and detailed analysis on our collected data are provided in the Appendix.

While our training data collection primarily targets iPhone screens, we assemble test sets for both iPhone and Android platforms. For each task, we select 25 test screens from iPhone and 5 from Android. Due to overlaps in images across different tasks, the total number of unique images amounts to 56 for iPhone and 13 for Android. For evaluation, we randomly select 2 QA pairs for the conversational tasks, creating two distinct test instances with precisely one question in each input. Utilizing these test images, we formulate 20/40/38/20 questions iPhone and 5/10/10/10 questions for the Android sets, for the four tasks respectively.

## 4 Experiments

**Model Training.** During training, both the decoder and the projection layer are updated while the vision encoder is kept frozen. All the training data is formatted into the instruction-following format, and the training objective is the same as in Ferret. In total, our training mixture has 250K samples. It takes 26 hours to train on 8 A100 GPUs.

	Public Benchmark			Elementary Tasks				Advanced Tasks	
	S2W	WiC	TaP	Ref-i	Ref-A	Grd-i	Grd-A	iPhone	Android
Spotlight [27]	106.7	141.8	<b>88.4</b>	-	-	-	-	-	-
Ferret [46]	17.6	1.2	46.2	13.3	13.9	8.6	12.9	20.0	20.7
Ferret-UI	<b>113.4</b>	<b>142.0</b>	78.4	<b>80.5</b>	<b>82.4</b>	<b>79.4</b>	<b>83.5</b>	73.4	80.5
GPT-4V [1]	34.8	23.5	47.6	61.3	37.7	70.3	4.7	<b>114.3</b>	<b>128.2</b>

**Table 2:** Results of Ferret-UI and baseline models. *S2W*: screen2words, *WiC*: widget captions, *TaP*: taperception. We report the CIDEr score for S2W and WiC and F1 for TaP. For elementary and advanced tasks, we report the averaged performance of corresponding tasks. “i”: iPhone, “A”: Android, “Ref”: Referring, “Grd”: Grounding.

## 4.1 Results

We compare the performance of Ferret-UI, Ferret<sup>2</sup>, and GPT-4V for all tasks. We also include Fuyu [5] and CogAgent’s [18] performances on advanced tasks.<sup>3</sup> Results are summarized in Tab. 2, where the average performance within a category is reported. Performance breakdown for elementary and advanced tasks is shown in Fig. 5a and Tab. 5b, respectively.

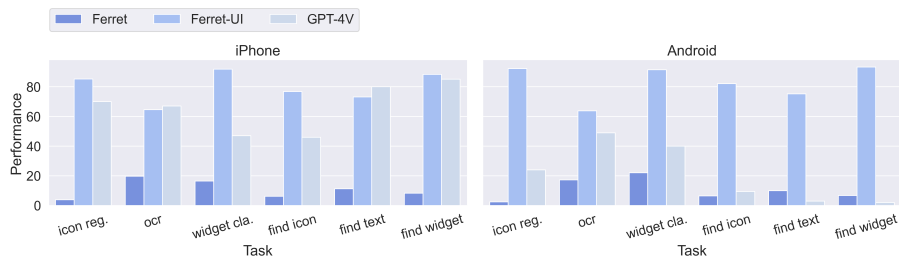
**Public Benchmark from Spotlight [27].** Compared to Spotlight, Ferret-UI demonstrates superior performance in *S2W* and *WiC*, even though Spotlight uses 80 million web page screenshots and 2.69 million mobile screenshots for pre-training. Ferret-UI performance falls short on *TaP* but still competitive; our studies further suggests that this could be due to the noisiness of the taperception labels. Detailed analysis is provided in the Appendix.

**Results on Elementary UI Tasks.** The average performance of all referring and grounding tasks is summarized in Tab. 2, and the performance breakdown for each task is shown in Fig. 5a. For referring tasks, we report exact match accuracy for OCR and accuracy for icon recognition and widget classification. For each grounding task, we also report the accuracy, where a correct bounding box is one that has an Intersection over Union (IoU) with the label greater than the threshold (0.5). Widget listing performances are not included in the average as we treat it as an auxiliary task.

Ferret-UI outperforms Ferret and GPT-4V in most elementary tasks except for iPhone *OCR* and *find text*, as GPT-4V is likely to be exposed to more text data. While GPT-4V demonstrates decent performance on iPhone tasks, its performances on Android tasks, especially grounding tasks, are significantly worse. Examining the predictions shows that Android screens have more numerous and smaller widgets, making the grounding tasks more challenging. Furthermore,

<sup>2</sup> For Ferret, we include the pre-defined classes for icon classification and widget classification in the prompts while the remaining prompts are the same as Ferret-UI.

<sup>3</sup> For GPT-4V, we sample a random subset of 100 instances for the Spotlight and elementary tasks for cost efficiency. For GPT-4V evaluation, we follow [43] by overlaying indexed bounding boxes of UI elements as visual prompts. Consequently, in grounding tasks, GPT-4V is enabled to make selections from among these candidate boxes. We detail the effort in Appendix.



(a) Elementary task performance comparison.

	iPhone					Android				
	DetDes	ConvP	ConvI	FuncIn	Avg	DetDes	ConvP	ConvI	FuncIn	Avg
Ferret [46]	2.5	34.7	23.7	19.1	20.0	2.0	33.9	24.9	21.9	20.7
Fuyu [5]	5.0	24.6	18.8	35.7	21.0	2.0	20.8	44.5	36.1	25.9
CogAgent [18]	53.1	59.7	74.8	71.9	64.9	28.0	58.5	90.1	<b>90.5</b>	66.8
Ferret-UI	64.5	75.0	77.5	76.5	73.4	90.8	72.8	79.3	79.2	80.5
GPT-4V [1]	<b>66.8</b>	<b>105.6</b>	<b>198.5</b>	<b>86.3</b>	<b>114.3</b>	<b>126.6</b>	<b>109.4</b>	<b>188.6</b>	88.3	<b>128.2</b>

(b) Advanced task performance comparison.

**Fig. 5:** Task performance breakdown. *DetDes*: detailed description, *ConvP*: conversation perception, *ConvI*: conversation interaction, *FuncIn*: function inference. Numerous small widgets present on the Android screen make it more challenging for referring and grounding, while Ferret-UI continues to outperform other models in the majority of elementary tasks and maintains parity with GPT-4V.

Ferret-UI’s zero-shot performance on the Referring Expression Comprehension task from UIBert [4] is 76% when we frame it as the *find widget* task.

**Results on Advanced Tasks.** The breakdown of task performance for advanced tasks is shown in Tab. 5b. As the advanced tasks require open-ended responses, we use GPT-4 to score both the label and the prediction. We report *score for prediction* over *score for label* as a percentage.

Ferret-UI exhibits commendable performance on advanced tasks for both platforms, despite the absence of Android-specific data in its training dataset. This suggests a notable transferability of UI knowledge across different operating systems, with Ferret-UI outperforming the average. While Fuyu tends to generate answers that are generally relevant, its responses lack the detail and precision exhibited by Ferret-UI. Conversely, GPT-4V secures higher scores across all tasks by consistently delivering more detailed responses than Ferret-UI, a characteristic that aligns with the preferences of the model evaluator (GPT-4).

## 4.2 Ablation Studies

**Advanced Tasks.** The design motivation behind elementary tasks is to enhance the model’s visual and spatial understanding of basic UI elements. Our hypothesis is that this enhanced understanding can help with carrying out advanced tasks. We investigate this hypothesis by studying whether elementary tasks affect the model’s performances on advanced tasks and we present the re-

	iPhone	Android		S2W	WiC	TaP
Adv. task only	64.6	64.3	Spotlight [27]	106.7	141.8	<b>88.4</b>
+ iPhone elem.	70.3	68.6	Balanced TaP labels	111.7	133.8	76.5
+ Android elem.	70.2	75.3	Spotlight tasks only	111.3	138.7	77.6
+ both as in 2	<b>73.4</b>	<b>80.5</b>	+ Android elem. tasks	111.3	138.0	76.8
			+ iPhone elem. tasks	112.4	138.9	74.8
			+ both	111.3	138.7	76.0
			Full mixture from 2	<b>113.4</b>	<b>142.0</b>	78.4

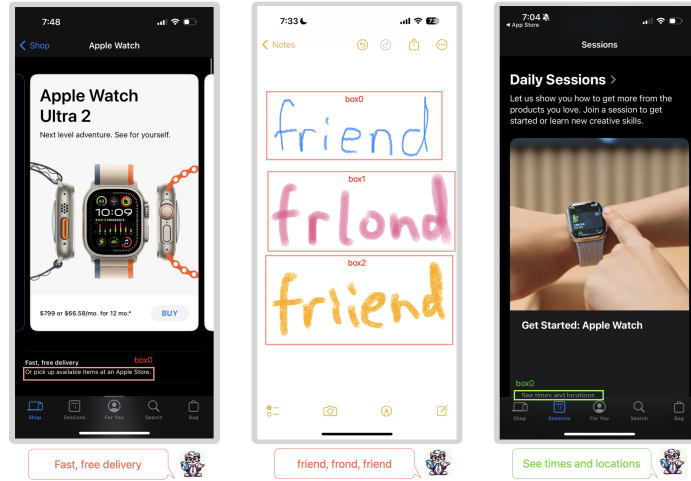
(a) **Advanced Tasks Ablation.** Performance on iPhone and Android advanced tasks. The training configurations are mixing advanced task with no other data, with iPhone elementary tasks only, Android elementary tasks only, or both.

(b) **Spotlight Tasks Ablation.** Performance on *S2W*, *WiC*, *TaP* tasks. For the balanced TaP labels experiment, we up-sample the minority class.

**Table 3:** Ablation studies on the factors that impact (a) advanced tasks performance and (b) Spotlight tasks performance, respectively.

sults in Tab. 3a. We see that with only advanced task data, the performance is 64% for both platforms. The performance of advanced tasks on iPhone shows a consistent improvement of 5% with the addition of either iPhone or Android elementary tasks. Conversely, when adding iPhone elementary tasks, Android advanced task performance improves by approximately 4%, while the addition of Android elementary tasks leads to a 9% improvement. With both sets of elementary tasks, iPhone and Android advanced tasks see a further improvement of 3% and 5% respectively above the best single set added. These observations support our hypothesis that elementary tasks provide the model with enhanced visual and spatial understanding that facilitates advanced tasks.

**Spotlight Tasks.** Motivated by a desire to explore the impact of different data configurations on Spotlight task performance, we specifically investigate whether adding elementary tasks data could enhance the model performance, given that these tasks are designed to improve the visual and spatial comprehension of screens. As shown in Tab. 3b, the addition of elementary task data—whether exclusively from Android, iPhone, or a combination of both—does not significantly alter performance across the three Spotlight tasks. This may be attributed to the short and highly specialized UI-centric vocabulary used in responses to elementary tasks, contrasting with the response style demanded by Spotlight tasks. Optimal results for Spotlight tasks were observed when data from advanced tasks were integrated alongside all elementary tasks, even though the advanced task data was exclusively derived from iPhone screens. Notably, this configuration yields a 4-point CiDER score improvement in widget captions upon the addition of advanced task data. We postulate that the free-response format of advanced task answers, which necessitates a more sophisticated set of skills for execution, aligns more closely with the requirements of Spotlight tasks. These tasks demand a comprehensive understanding beyond that of recognizing individual UI elements, as is common in elementary tasks. Moreover, executing advanced tasks requires more sophisticated skills than understanding one specific UI element on the screen as in elementary tasks. Therefore, it is expected that these skills help with Spotlight tasks, which are more complex than ele-



**Fig. 6: OCR Analyses.** *Left:* predict nearby text instead of targeted region. *Middle:* a tendency to predict valid words. *Right:* Ferret-UI correctly reads cut-off text, while detection model produces wrong label.

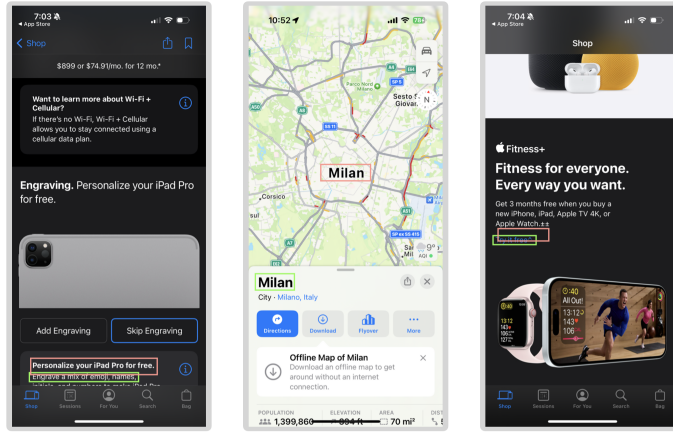
mentary tasks but less convoluted than advanced tasks. In conclusion, the task formulation plays a bigger role than the platform of the data that is added.

### 4.3 Result Analysis: Elementary UI Tasks

**Referring Tasks Analyses.** In analyzing Ferret-UI’s capabilities, we specifically focus on OCR predictions. The OCR analysis reveals three notable observations, as depicted in Fig. 6. Firstly, the model predicts a neighboring text instead of the text in the targeted region. This is common for smaller texts and texts very close to other texts. Secondly, the model exhibits a tendency to predict actual words rather than merely deciphering characters displayed on the screen. This observation is in line with the semantic-reliance observation of LLMs made in some existing work [32]. On UI screens, phonetically crafted words that are commonly used as brand titles largely fall under this category. Thirdly, Ferret-UI demonstrates the ability to accurately predict text that is partially cut-off, even in instances where the OCR model returns incorrect texts.

**Grounding Tasks Analyses.** Using *find text* predictions, as depicted in Fig. 7, we further elucidate observations from grounding tasks. Echoing the initial observation from the *OCR* analysis, the model may erroneously highlight a piece of text adjacent to the targeted area. Additionally, the occurrence of multiple instances of identical texts suggests the potential for expanding future methods to encompass a range of answers from a singular box to multiple boxes, thereby enhancing the model’s utility and accuracy in complex text-finding scenarios.





**Fig. 7: Find Text Analyses.** *Left:* a neighboring text is mis-identified as the target. *Middle:* multiple occurrences of the same text. *Right:* predicted boxes not precise.

#### 4.4 Result Analysis: Advanced UI Tasks

**Grounded conversation.** Engaging in grounded conversation is Ferret’s unique capability. To better understand the quality of the output bounding boxes in terms of correctness and relevance, we manually grade all output boxes in both Ferret-UI and GPT-4V’s *conversation interaction* outputs. The accuracies for Ferret-UI and GPT-4V are 91.7% and 93.4% respectively. Considering Ferret-UI is generating raw coordinates whereas GPT-4V chooses from a set of pre-defined boxes, Ferret-UI’s grounding ability on UI screens is noteworthy. Even though Ferret-UI’s received evaluation score falls short to GPT-4V, from inspecting the predictions as in Fig. 8, we notice that GPT-4V tends to provide extra information that may not be relevant to the question. However, these detailed answers are more favored in scoring than Ferret-UI’s concise answers.

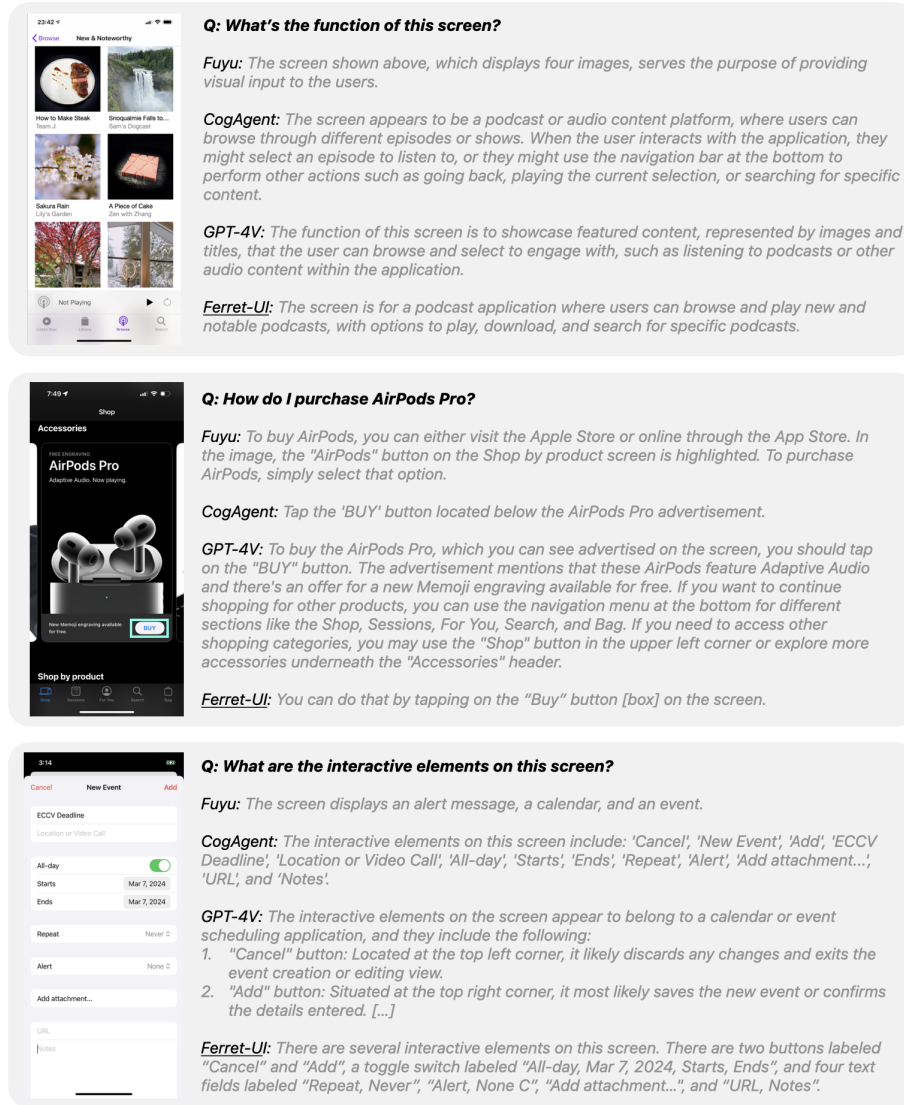
**UI detection model is a bottleneck.** Given that both our elementary and advanced tasks are predicated upon detections of UI elements, Ferret-UI is not able to learn aspects of screens that are not detected, such as colors, design, usability, and UI elements that the detection model misses (eg, topmost time, wifi, battery, etc). For example, in generating detailed descriptions, GPT-4V is capable of noting “The overall design conforms to Apple’s aesthetic with a minimalistic, clean, dark theme.” a level of insight Ferret-UI is not trained to offer due to its reliance on detected elements alone.

**Set-of-Mark (SoM) Prompting of GPT-4V.** In our evaluation of GPT-4V, we employ the Set-of-Mark (SoM) prompting technique, which encounters several notable limitations. Firstly, its effectiveness diminishes in scenarios involving a multitude of small UI elements, a common occurrence in Android detection tasks. The small size of some UI components means that the addition of labels may obscure original content or even extend beyond the intended areas. Secondly,

constraining the evaluation to a pre-defined set of candidate regions restricts the model’s ability to reference any given region freely. In the middle example shown in Fig. 8, the UI detection model treats the entire middle section as one element, covering the texts, image, and the Buy button. Therefore, the model is not able to refer to the "BUY" button on its own in its responses, since it is considered part of a collective detection group.

## 5 Conclusion

In summary, we introduce Ferret-UI, a specialized MLLM designed to enhance comprehension and interaction with mobile UI screens. Through careful curation of training samples encompassing a diverse range of basic and advanced UI tasks, Ferret-UI demonstrates remarkable proficiency in referring, grounding, and reasoning. The advent of these enhanced capabilities promises substantial advancements for a multitude of downstream UI applications, thereby amplifying the potential benefits afforded by Ferret-UI in this domain.



**Fig. 8:** Visualization results of advanced tasks (top to bottom: *function inference*, *conversation interaction*, *conversation perception*) to illustrate the differences among various models (Fuyu v.s. CogAgent v.s. GPT-4V v.s. Ferret-UI.)

## References

1. Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F.L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al.: Gpt-4 technical report. arXiv preprint arXiv:2303.08774 (2023)
2. Amalfitano, D., Fasolino, A.R., Tramontana, P.: A gui crawling-based technique for android mobile application testing. In: 2011 IEEE fourth international conference on software testing, verification and validation workshops. pp. 252–261. IEEE (2011)
3. Anil, R., Dai, A.M., Firat, O., Johnson, M., Lepikhin, D., Passos, A., Shakeri, S., Taropa, E., Bailey, P., Chen, Z., et al.: Palm 2 technical report. arXiv preprint arXiv:2305.10403 (2023)
4. Bai, C., Zang, X., Xu, Y., Sunkara, S., Rastogi, A., Chen, J., y Arcas, B.A.: Uibert: Learning generic multimodal representations for ui understanding (2021)
5. Bavishi, R., Elsen, E., Hawthorne, C., Nye, M., Odena, A., Somani, A., Taşlılar, S.: Introducing our multimodal models (2023), <https://www.adept.ai/blog/fuyu-8b>
6. Burns, A., Arsan, D., Agrawal, S., Kumar, R., Saenko, K., Plummer, B.A.: A dataset for interactive vision-language navigation with unknown command feasibility. In: European Conference on Computer Vision. pp. 312–328. Springer (2022)
7. Chen, K., Zhang, Z., Zeng, W., Zhang, R., Zhu, F., Zhao, R.: Shikra: Unleashing multimodal llm’s referential dialogue magic. arXiv preprint arXiv:2306.15195 (2023)
8. Cheng, K., Sun, Q., Chu, Y., Xu, F., Li, Y., Zhang, J., Wu, Z.: Seeclick: Harnessing gui grounding for advanced visual gui agents (2024)
9. Dai, W., Li, J., Li, D., Tiong, A.M.H., Zhao, J., Wang, W., Li, B., Fung, P., Hoi, S.: Instructblip: Towards general-purpose vision-language models with instruction tuning. arXiv preprint arXiv:2305.06500 (2023)
10. Deka, B., Huang, Z., Franzen, C., Hirschman, J., Afegan, D., Li, Y., Nichols, J., Kumar, R.: Rico: A mobile app dataset for building data-driven design applications. In: Proceedings of the 30th annual ACM symposium on user interface software and technology. pp. 845–854 (2017)
11. Deng, X., Gu, Y., Zheng, B., Chen, S., Stevens, S., Wang, B., Sun, H., Su, Y.: Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing Systems* **36** (2024)
12. Driess, D., Xia, F., Sajjadi, M.S., Lynch, C., Chowdhery, A., Ichter, B., Wahid, A., Tompson, J., Vuong, Q., Yu, T., et al.: Palm-e: An embodied multimodal language model. arXiv preprint arXiv:2303.03378 (2023)
13. Edwards, W.K., Mynatt, E.D., Stockton, K.: Access to graphical interfaces for blind users. *Interactions* **2**(1), 54–67 (1995)
14. Gao, D., Ji, L., Bai, Z., Ouyang, M., Li, P., Mao, D., Wu, Q., Zhang, W., Wang, P., Guo, X., et al.: Assistgui: Task-oriented desktop graphical user interface automation. arXiv preprint arXiv:2312.13108 (2023)
15. Gu, A., Dao, T.: Mamba: Linear-time sequence modeling with selective state spaces. arXiv preprint arXiv:2312.00752 (2023)
16. Gur, I., Furuta, H., Huang, A., Safdari, M., Matsuo, Y., Eck, D., Faust, A.: A real-world webagent with planning, long context understanding, and program synthesis. arXiv preprint arXiv:2307.12856 (2023)
17. Gur, I., Rueckert, U., Faust, A., Hakkani-Tur, D.: Learning to navigate the web. arXiv preprint arXiv:1812.09195 (2018)

18. 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. arXiv preprint arXiv:2312.08914 (2023)
19. Huang, S., Dong, L., Wang, W., Hao, Y., Singhal, S., Ma, S., Lv, T., Cui, L., Mohammed, O.K., Patra, B., Liu, Q., Aggarwal, K., Chi, Z., Bjorck, J., Chaudhary, V., Som, S., Song, X., Wei, F.: Language is not all you need: Aligning perception with language models (2023)
20. Jiang, A.Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D.S., Casas, D.d.l., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., et al.: Mistral 7b. arXiv preprint arXiv:2310.06825 (2023)
21. Jiang, Y., Schoop, E., Swearngin, A., Nichols, J.: Iluvui: Instruction-tuned language-vision modeling of uis from machine conversations (2023)
22. Jiang, Z., Kuang, R., Gong, J., Yin, H., Lyu, Y., Zhang, X.: What makes a great mobile app? a quantitative study using a new mobile crawler. In: 2018 IEEE Symposium on Service-Oriented System Engineering (SOSE). pp. 222–227. IEEE (2018)
23. Kim, G., Baldi, P., McAleer, S.: Language models can solve computer tasks (2023)
24. Lai, X., Tian, Z., Chen, Y., Li, Y., Yuan, Y., Liu, S., Jia, J.: Lisa: Reasoning segmentation via large language model. arXiv preprint arXiv:2308.00692 (2023)
25. Lee, K., Joshi, M., Turc, I.R., Hu, H., Liu, F., Eisenschlos, J.M., Khandelwal, U., Shaw, P., Chang, M.W., Toutanova, K.: Pix2struct: Screenshot parsing as pretraining for visual language understanding. In: ICML (2023)
26. Li, B., Zhang, Y., Chen, L., Wang, J., Yang, J., Liu, Z.: Otter: A multi-modal model with in-context instruction tuning. arXiv preprint arXiv:2305.03726 (2023)
27. Li, G., Li, Y.: Spotlight: Mobile ui understanding using vision-language models with a focus (2023)
28. Li, Y., He, J., Zhou, X., Zhang, Y., Baldridge, J.: Mapping natural language instructions to mobile ui action sequences. arXiv preprint arXiv:2005.03776 (2020)
29. Linares-Vásquez, M., Moran, K., Poshyvanyk, D.: Continuous, evolutionary and large-scale: A new perspective for automated mobile app testing. In: 2017 IEEE International Conference on Software Maintenance and Evolution (ICSME). pp. 399–410. IEEE (2017)
30. Liu, E.Z., Guu, K., Pasupat, P., Shi, T., Liang, P.: Reinforcement learning on web interfaces using workflow-guided exploration. arXiv preprint arXiv:1802.08802 (2018)
31. Liu, H., Li, C., Wu, Q., Lee, Y.J.: Visual instruction tuning (2023)
32. Liu, Y., Li, Z., Yang, B., Li, C., Yin, X., Lin Liu, C., Jin, L., Bai, X.: On the hidden mystery of ocr in large multimodal models (2024)
33. OpenAI, :, Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F.L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., Avila, R., Babuschkin, I., Balaji, S., Balcom, V., Baltescu, P., Bao, H., Bavarian, M., Belgum, J., Bello, I., Berdine, J., Bernadett-Shapiro, G., Berner, C., Bogdonoff, L., Boiko, O., Boyd, M., Brakman, A.L., Brockman, G., Brooks, T., Brundage, M., Button, K., Cai, T., Campbell, R., Cann, A., Carey, B., Carlson, C., Carmichael, R., Chan, B., Chang, C., Chantzis, F., Chen, D., Chen, S., Chen, R., Chen, J., Chen, M., Chess, B., Cho, C., Chu, C., Chung, H.W., Cummings, D., Currier, J., Dai, Y., Decareaux, C., Degry, T., Deutsch, N., Deville, D., Dhar, A., Dohan, D., Dowling, S., Dunning, S., Ecoffet, A., Eleti, A., Eloundou, T., Farhi, D., Fedus, L., Felix, N., Fishman, S.P., Forte, J., Fulford, I., Gao, L., Georges, E., Gibson, C., Goel, V., Gogineni, T., Goh, G., Gontijo-Lopes, R., Gordon, J., Grafstein, M., Gray, S., Greene, R.,

- Gross, J., Gu, S.S., Guo, Y., Hallacy, C., Han, J., Harris, J., He, Y., Heaton, M., Heidecke, J., Hesse, C., Hickey, A., Hickey, W., Hoeschele, P., Houghton, B., Hsu, K., Hu, S., Hu, X., Huizinga, J., Jain, S., Jain, S., Jang, J., Jiang, A., Jiang, R., Jin, H., Jin, D., Jomoto, S., Jonn, B., Jun, H., Kaftan, T., Łukasz Kaiser, Kamali, A., Kanitscheider, I., Keskar, N.S., Khan, T., Kilpatrick, L., Kim, J.W., Kim, C., Kim, Y., Kirchner, J.H., Kiros, J., Knight, M., Kokotajlo, D., Łukasz Kondraciuk, Kondrich, A., Konstantinidis, A., Kosic, K., Krueger, G., Kuo, V., Lampe, M., Lan, I., Lee, T., Leike, J., Leung, J., Levy, D., Li, C.M., Lim, R., Lin, M., Lin, S., Litwin, M., Lopez, T., Lowe, R., Lue, P., Makanju, A., Malfacini, K., Manning, S., Markov, T., Markovski, Y., Martin, B., Mayer, K., Mayne, A., McGrew, B., McKinney, S.M., McLeavey, C., McMillan, P., McNeil, J., Medina, D., Mehta, A., Menick, J., Metz, L., Mishchenko, A., Mishkin, P., Monaco, V., Morikawa, E., Mossing, D., Mu, T., Murati, M., Murk, O., Mély, D., Nair, A., Nakano, R., Nayak, R., Neelakantan, A., Ngo, R., Noh, H., Ouyang, L., O’Keefe, C., Pachocki, J., Paino, A., Palermo, J., Pantuliano, A., Parascandolo, G., Parish, J., Parparita, E., Passos, A., Pavlov, M., Peng, A., Perelman, A., de Avila Belbute Peres, F., Petrov, M., de Oliveira Pinto, H.P., Michael, Pokorny, Pocrass, M., Pong, V.H., Powell, T., Power, A., Power, B., Proehl, E., Puri, R., Radford, A., Rae, J., Ramesh, A., Raymond, C., Real, F., Rimbach, K., Ross, C., Rotsted, B., Roussez, H., Ryder, N., Saltarelli, M., Sanders, T., Santurkar, S., Sastry, G., Schmidt, H., Schnurr, D., Schulman, J., Selsam, D., Sheppard, K., Sherbakov, T., Shieh, J., Shoker, S., Shyam, P., Sidor, S., Sigler, E., Simens, M., Sitkin, J., Slama, K., Sohl, I., Sokolowsky, B., Song, Y., Staudacher, N., Such, F.P., Summers, N., Sutskever, I., Tang, J., Tezak, N., Thompson, M.B., Tillet, P., Tootoonchian, A., Tseng, E., Tuggle, P., Turley, N., Tworek, J., Uribe, J.F.C., Vallone, A., Vijayvergiya, A., Voss, C., Wainwright, C., Wang, J.J., Wang, A., Wang, B., Ward, J., Wei, J., Weinmann, C., Welihinda, A., Welinder, P., Weng, J., Weng, L., Wiethoff, M., Willner, D., Winter, C., Wolrich, S., Wong, H., Workman, L., Wu, S., Wu, J., Wu, M., Xiao, K., Xu, T., Yoo, S., Yu, K., Yuan, Q., Zaremba, W., Zellers, R., Zhang, C., Zhang, M., Zhao, S., Zheng, T., Zhuang, J., Zhuk, W., Zoph, B.: Gpt-4 technical report (2024)
34. Peng, Z., Wang, W., Dong, L., Hao, Y., Huang, S., Ma, S., Wei, F.: Kosmos-2: Grounding multimodal large language models to the world. arXiv preprint arXiv:2306.14824 (2023)
  35. 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)
  36. Shaw, P., Joshi, M., Cohan, J., Berant, J., Pasupat, P., Hu, H., Khandelwal, U., Lee, K., Toutanova, K.N.: From pixels to ui actions: Learning to follow instructions via graphical user interfaces. *Advances in Neural Information Processing Systems* **36** (2024)
  37. Shi, T., Karpathy, A., Fan, L., Hernandez, J., Liang, P.: World of bits: An open-domain platform for web-based agents. In: International Conference on Machine Learning. pp. 3135–3144. PMLR (2017)
  38. Sun, Q., Yu, Q., Cui, Y., Zhang, F., Zhang, X., Wang, Y., Gao, H., Liu, J., Huang, T., Wang, X.: Generative pretraining in multimodality. arXiv preprint arXiv:2307.05222 (2023)
  39. Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al.: Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971 (2023)

40. Wang, J., Xu, H., Ye, J., Yan, M., Shen, W., Zhang, J., Huang, F., Sang, J.: Mobile-agent: Autonomous multi-modal mobile device agent with visual perception (2024)
41. 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)
42. Yan, A., Yang, Z., Zhu, W., Lin, K., Li, L., Wang, J., Yang, J., Zhong, Y., McAuley, J., Gao, J., et al.: Gpt-4v in wonderland: Large multimodal models for zero-shot smartphone gui navigation. arXiv preprint arXiv:2311.07562 (2023)
43. Yang, J., Zhang, H., Li, F., Zou, X., Li, C., Gao, J.: Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v. arXiv preprint arXiv:2310.11441 (2023)
44. Yang, Z., Li, L., Lin, K., Wang, J., Lin, C.C., Liu, Z., Wang, L.: The dawn of llms: Preliminary explorations with gpt-4v(ision) (2023)
45. Ye, Q., Xu, H., Xu, G., Ye, J., Yan, M., Zhou, Y., Wang, J., Hu, A., Shi, P., Shi, Y., et al.: mplug-owl: Modularization empowers large language models with multimodality. arXiv preprint arXiv:2304.14178 (2023)
46. You, H., Zhang, H., Gan, Z., Du, X., Zhang, B., Wang, Z., Cao, L., Chang, S.F., Yang, Y.: Ferret: Refer and ground anything anywhere at any granularity (2023)
47. Zhan, Z., Zhang, A.: You only look at screens: Multimodal chain-of-action agents. arXiv preprint arXiv:2309.11436 (2023)
48. Zhang, C., Yang, Z., Liu, J., Han, Y., Chen, X., Huang, Z., Fu, B., Yu, G.: Appagent: Multimodal agents as smartphone users (2023)
49. Zhang, S., Sun, P., Chen, S., Xiao, M., Shao, W., Zhang, W., Chen, K., Luo, P.: Gpt4roi: Instruction tuning large language model on region-of-interest. arXiv preprint arXiv:2307.03601 (2023)
50. Zhang, X., de Greef, L., Swearngin, A., White, S., Murray, K., Yu, L., Shan, Q., Nichols, J., Wu, J., Fleizach, C., Everitt, A., Bigham, J.P.: Screen recognition: Creating accessibility metadata for mobile applications from pixels (2021)
51. Zhao, Y., Lin, Z., Zhou, D., Huang, Z., Feng, J., Kang, B.: Bubogpt: Enabling visual grounding in multi-modal llms. arXiv preprint arXiv:2307.08581 (2023)
52. Zheng, B., Gou, B., Kil, J., Sun, H., Su, Y.: Gpt-4v (ision) is a generalist web agent, if grounded. arXiv preprint arXiv:2401.01614 (2024)
53. Zheng, L., Chiang, W.L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E.P., Zhang, H., Gonzalez, J.E., Stoica, I.: Judging llm-as-a-judge with mt-bench and chatbot arena (2023)
54. Zheng, L., Wang, R., Wang, X., An, B.: Synapse: Trajectory-as-exemplar prompting with memory for computer control (2024)
55. Zhu, D., Chen, J., Shen, X., Li, X., Elhoseiny, M.: Minigpt-4: Enhancing vision-language understanding with advanced large language models. arXiv preprint arXiv:2304.10592 (2023)

## A Elementary Task Data Generation Details

Additional details in elementary task data generation are as follows:

- In our data generation process, we merge the two distinct classes—‘Checked’ and ‘Unchecked’—found in the original detection labels for both *Checkboxes* and *Toggles*.
- For widget listing, the answer starts with a common phrase: *UI widgets present in this screen include*. Each element is formatted as “{displayed text} {UI type}” (e.g., “login button”), except for text elements, which are formatted as “Text displaying {displayed text}”.
- For OCR, we consider text with fewer than 10 tokens. If the text is exactly one token, the length needs to be 2 or greater to be included.
- For tasks such as *find text*, *find icons*, and *find widget*, it is common to encounter screens containing multiple instances of the same UI element (e.g., multiple login buttons). We employ a filtering mechanism that excludes samples involving UI elements with multiple occurrences within a single screen.
- The size of the test set is determined by selecting the smaller value between 5k and the total number of generated test instances.

## B Advanced Task Data Quality Analysis

We conduct a thorough analysis of the quality of our collected data for advanced tasks and provide comprehensive statistics. The vocabulary size for each task is as follows: 30,866 for *detailed description*, 15,666 for *conversation perception*, 12,092 for *conversation interaction*, and 14,896 for *function inference*.

In the realm of *conversation interaction*, we observe 33,649 question turns and 29,933 answer turns. Among these, 15 question turns include bounding boxes, whereas all answer turns include bounding boxes. We compile the most frequently occurring tri-grams for questions and answers in both conversation tasks. Notably, in *conversation perception* questions, the top tri-grams include phrases like *are there any*”, *where is the*”, and *what is the*”, while those for interactions comprise phrases like *How can I*”, *I want to*”, and *Can I do*”. Similarly, in perception answers, prominent tri-grams consist of expressions such as *bottom of the*”, *at the top*”, and *there is a*”, while interaction answers primarily feature tri-grams like *by tapping on*”, *tapping on the*”, and *can tap on*”.

We present detailed distributions of tri-grams in conversation data questions and answers in Fig. 9. This observation is consistent with our intended objectives for each conversation category, with perception focusing on visual elements and interaction emphasizing actions. Notably, from the interaction conversation answers, we observe that *tap* emerges as the predominant action. In future work, we aim to explore interactions involving other actions, such as scrolling, long-clicking, and entering text. The inclusion of two conversation categories aims to diversify conversation topics, although a clear-cut distinction between the two is not always feasible, and overlap between the categories may occur.



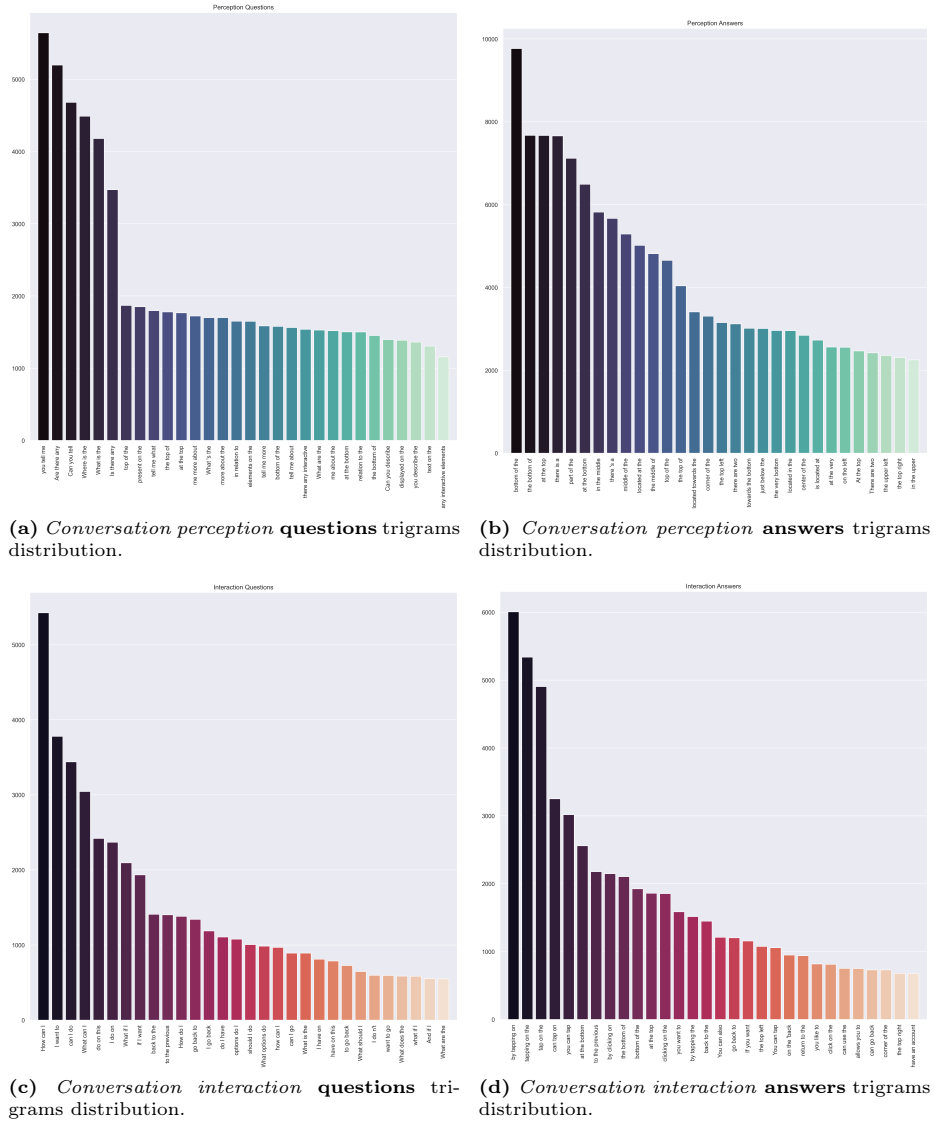
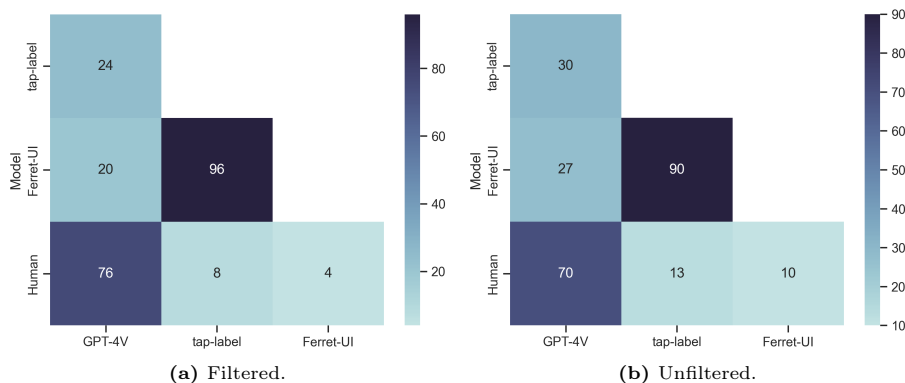


Fig. 9: Trigrams for collected conversation data questions and answers.

### C Taperception label analysis

We meticulously label 30 test samples for *taperception* and conduct a study on the correlation among our labels, *taperception* ground-truth labels, Ferret-UI outputs, and GPT-4V outputs. Among the 30 samples, 5 pose challenges in deciphering without direct interaction with the screen.

In Tab. 10, we present the percentage of agreement among different sources of predictions and labels. The term "filtered" denotes the set of 25 instances that are unambiguous, while "unfiltered" encompasses the entire 30 instances. We observe that our labels exhibit a high correlation with GPT-4V predictions, but differing significantly from the *taperception* dataset labels. This discrepancy underscores the complexity of predicting *tappability* solely based on single images, highlighting the inherent challenges in obtaining clear-cut labels for this task.



**Fig. 10:** Agreement between different sources of taperception predictions and labels. In unfiltered, we make the best educational guess for the one that are ambiguous. We observe that our human annotation correlate with GPT-4V (%76) far more than with taperception label (%8). Even though Ferret-UI performance on taperception falls behind compared to Spotlight, it could be due to the noisiness of labels.

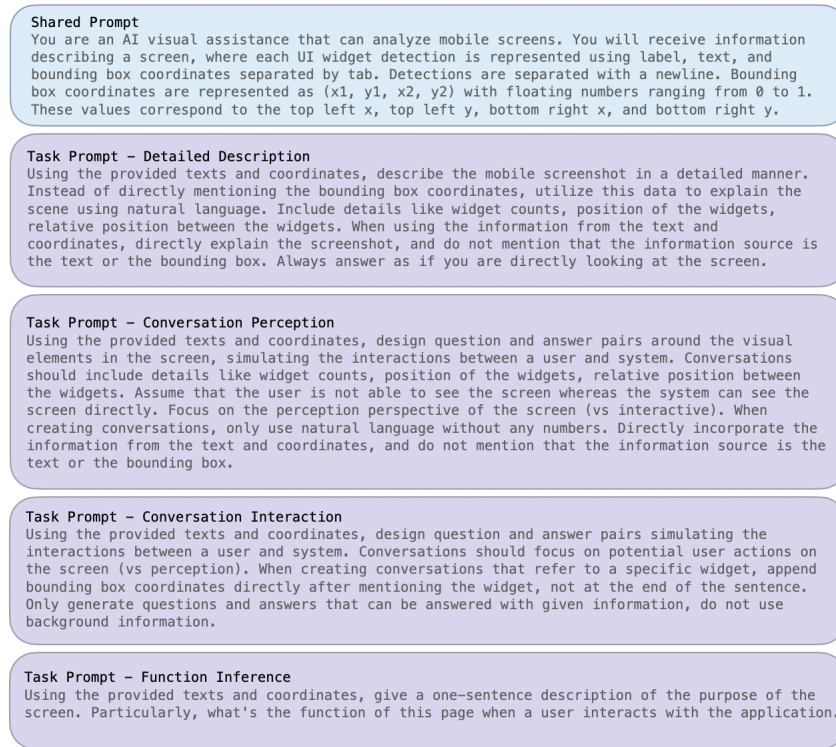
## D Advanced Task Generation Prompts

We present the prompts that we use to collect advanced task data from GPT-4 in Fig. D.

## E GPT-4V Evaluation Details

We detail the process of creating input for GPT-4V to tackle the UI tasks under scope.

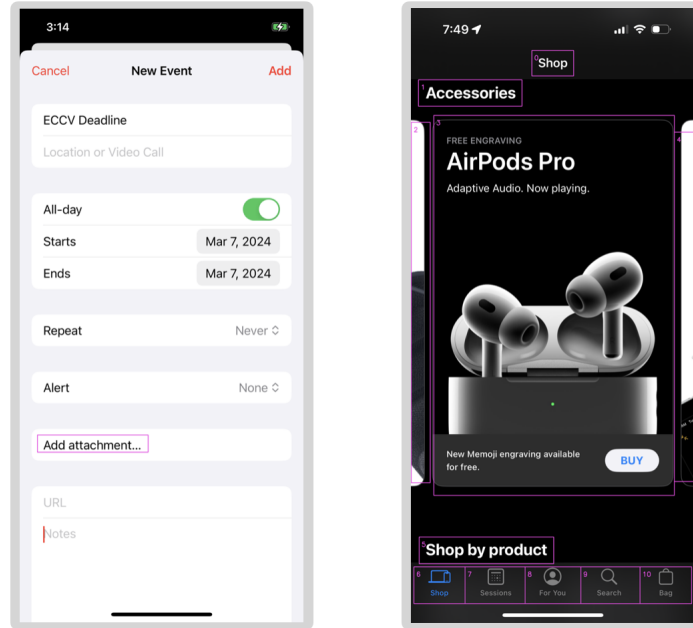
*[Input Images]* We first annotate the screenshots tailored to each specific task, ensuring that GPT-4V has sufficient contextual information to answer the questions. For tasks without any bounding boxes in input or output (*screen2words*, *widget captions*, and *Advanced Tasks*), we use the original images as the input.



**Fig. 11:** Prompts for GPT-4 in advanced task data generation.

For tasks that refer to **one** specific UI element using bounding box in the input, we put a magenta-colored bounding box on the image as the input, as shown in Fig. 12 left. For tasks that expect one or more bounding boxes in the output, our initial explorations confirm that GPT-4V is not able to provide bounding boxes in the output as it gives the answer, "*Unfortunately, I'm not able to provide the exact bounding box coordinates, as my capabilities are currently limited to describing images and discussing the content rather than interacting directly with the image to extract detailed metadata such as pixel coordinates.*") and proceed to answer the question in natural language. Therefore, for those tasks, we create an easier version where we ask GPT-4V to choose from a fixed set of candidates. Particularly, we follow Set-of-Mark prompting [43] where for each UI detection from our UI detection model, we use a magenta-colored bounding box to mark it in the screen and inside each box we assign a numeric label so that GPT-4V can refer to it. An example input image is shown in Fig. 12 right.

[Prompts] With the input images ready, we further modify the prompts to provide GPT-4V with all the necessary information to perform all the tasks successfully. For tap-perception, we instruct it to answer 'Yes.' or 'No.' only without

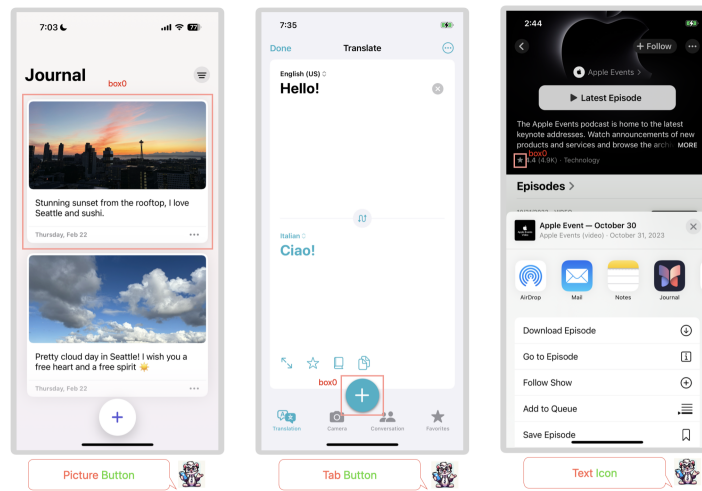


**Fig. 12:** GPT-4V input image examples. Left: used in referring task, where the question concerns on specific UI element. Right: used in grounding task, where GPT-4V refers to the UI elements by their assigned numeric labels.

any explanations. For widget captions, we instruct it to ‘*Answer in a few words.*’ For *icon recognition* and *widget classification*, we provide the list of all possible classes, and instruct it to output the class only without any explanations. For *OCR*, we instruct it to output the identified text only. For *find widget*, *find text*, *find icons*, we add to the prompt “*Candidate choices are indicated using magenta bounding boxes in the image and each box is associated with a numeric label. Output the numeric label as your answer, no need to explain.*”

## F Widget Classification Analysis

Similar to OCR analysis, we show three interesting observations in Fig 13. First, the model struggles when it needs to understand relationships among widgets. For example, if a large button is made up of a few sub-elements including Picture, Icon, and text, the model cannot see it as a unified widget but tend to predict it as the sub-element that occupies the largest space. In line of first observation, when a Tab or an Icon is seated on top of a row of tabs, it is highly likely be considered part of the tabs. Finally, we discover a common case where small Icons surrounded by texts are likely to be predicted as Text, this is consistent with the observation that small texts tend to be predicted as neighboring texts.



**Fig. 13: Widget Classification Analyses.** *Left:* a large Button consists of Picture, Icon, and Text mis-classified as a Picture. *Middle:* a button seated on top of a row of Tabs mis-classified as a Tab. *Right:* a small, text-surrounded icon being classified as text.

## G More Example Outputs

**Q:** Classify the type of UI widget located in the element [box0].  
**A:** Toggle

**Q:** What text is displayed on the widget [box1].  
**A:** Add attachment...

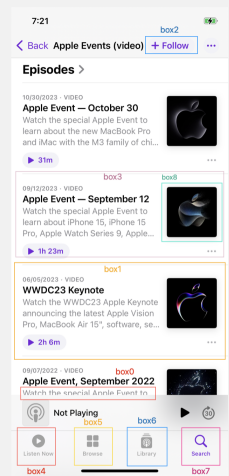
**Q:** How should I add a meeting link to this reminder?  
**A:** You can add a meeting link by tapping on the "Add attachment..." button [box1]

**Q:** This is not an All-day event, how should I reflect this?  
**A:** You can change the event status by tapping on the "All-day, Mar 7, 2024, Starts, Ends, Repeat, Never" button [box5]. If it's not an all-day event, you can select "Never" in the "Repeat" field [box6].

**Q:** What are the interactive elements on this screen?  
**A:** There are several interactive elements on this screen. There are two buttons labeled "Cancel" and "Add", a toggle switch labeled "All-day, Mar 7, 2024, Starts, Ends, Repeat, Never", and four text fields labeled "Repeat, Never", "Alert, None C", "URL, Notes".

**Q:** Where is the Cancel button?  
**A:** Cancel button [box3]

**Q:** Where is the Add button?  
**A:** Add button [box4]



7:21

Apple Events (video) + Follow

Episodes >

Apple Event — October 30  
Watch the special Apple Event to learn about the new MacBook Pro and Mac with the M3 family of ch...

Apple Event — September 12  
Watch the special Apple Event to learn about iPhone 15, iPhone 15 Pro, Apple Watch Series 9, Apple...

WWDC23 Keynote  
Watch the WWDC23 Apple Keynote announcing the latest Apple Vision Pro, MacBook Air 15", software, se...

Apple Event, September 2022  
Watch the special Apple Event to learn about the latest Apple products and services.

Not Playing

Listen Now Browse Library Search

Q: What text is displayed on the widget [box0]?  
A: Apple Event, September 2020

Q: Classify the type of UI widget located in the element [box1].  
A: Button

Q: Classify the type of UI widget located in the element [box8].  
A: Picture

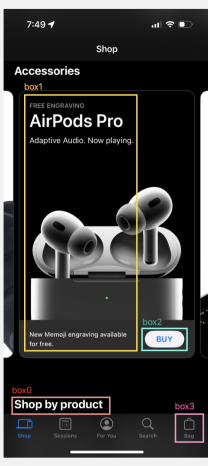
Q: How can I follow Apple Events?  
A: You can follow Apple Events by tapping on the "Follow" button [box2].

Q: What's the function of this screen when users interact with this app?  
A: The screen is for a video streaming application where users can view different episodes of Apple events, play them, and navigate through different tabs such as "Listen Now", "Browse", "Library", and "Search".

Q: How long is the September 12 Apple Event?  
A: The September 12 Apple Event is approximately 1 hour and 32 minutes long.

Q: How do I start watching the September 12 Apple Event?  
A: You can start watching the September 12 Apple Event by tapping on the corresponding button [box3].

Q: What tab options do I have on this screen?  
A: There are four tabs at the bottom of the screen. They are "Listen Now" [box4], "Browse" [box5], "Library" [box6], and "Search" [box7].



7:49

Shop

Accessories

FREE ENGRAVING  
AirPods Pro  
Adaptive Audio. Now playing.

Now Memoji engraving available for free.

BUY

Shop by product

Shop Store Mode Memoji Search Bag

Q: What's the function of this screen when the user interacts with application?  
A: The screen is for an online shopping application where users can browse and purchase products, specifically the AirPods Pro, with options to shop by product, accessories, and other features like Store Mode, Search, and Bag.

Q: Can I customize my AirPods?  
A: Yes, you can customize your AirPods by tapping on the "Shop by product" text [box0].

Q: Do I need to pay for the engraving?  
A: No, the text on the screen indicates that the engraving is free [box1]. However, you will need to pay for the AirPods Pro [box1].

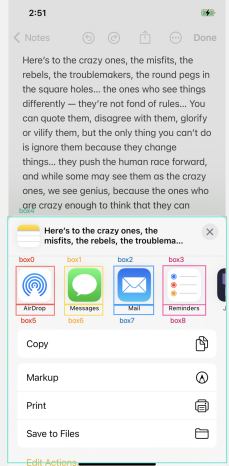
Q: Amazing! I want to buy AirPods!  
A: You can do that by tapping on the "Buy" button [box2] on the screen.

Q: Let me check what's in my bag.  
A: You can view the items in your bag by tapping on the "Bag" tab [box3].

Q: Please provide a detailed description for this screen.  
A: The screen appears to be a shopping page for a product, specifically the "AirPods Pro". The product is displayed prominently in the middle of the screen with a description stating "New Memoji engraving available, for free.". There is also a "Buy" button located towards the right side of the product description.

At the top of the screen, there is a "Shop" text centered, likely indicating the current page or section of the app. Below this, there is a section titled "Accessories", which is likely a category for the product.

At the bottom of the screen, there is a navigation bar with five tabs. From left to right, these are labeled "Shop", "Store Mode", "Memoji", "Search", and "Bag". The "Shop" tab is currently selected, as indicated by its position at the far left of the navigation bar.



2:51

Notes Done

Here's to the crazy ones, the misfits, the rebels, the troublemakers, the round pegs in the square holes... the ones who see things differently — they're not fond of rules... You can quote them, disagree with them, glorify or vilify them, but the only thing you can't do is ignore them because they change things... they push the human race forward, and while some may see them as the crazy ones, we see genius, because the ones who are crazy enough to think that they can

Here's to the crazy ones, the misfits, the rebels, the troublema...

AirDrop Messages Mail Reminders

Copy Markup Print Save to Files

Edit Actions

Q: Classify the type of icon displayed in the element [box0].  
A: wifi

Q: Classify the type of icon displayed in the element [box1].  
A: chat

Q: Classify the type of icon displayed in the element [box2].  
A: mail

Q: Classify the type of icon displayed in the element [box3].  
A: list

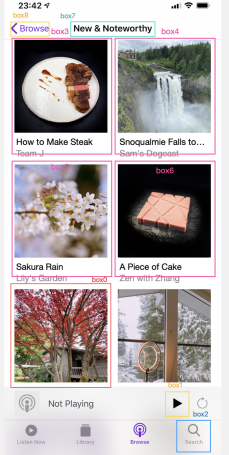
Q: Classify the type of UI widget located in the element [box4]?  
A: Dialog

Q: What text is displayed on the widget [box5]?  
A: AirDrop

Q: What text is displayed on the widget [box6]?  
A: Messages

Q: What text is displayed on the widget [box7]?  
A: Mail

Q: What text is displayed on the widget [box8]?  
A: Reminders



23:42

Browse New & Noteworthy

How to Make Steak Snoqualmie Falls to...

Sakura Rain A Piece of Cake

Not Playing

Listen Download Repeat Search

Q: Classify the type of UI widget located in the element [box0].  
A: Picture

Q: Where is the play icon located?  
A: play icon [box1]

Q: Where is the search tab located?  
A: search tab [box2]

Q: What podcasts are available on this page?  
A: There are several podcasts available on this page. They include "How to Make Steak" by Team J [box3], "Snowquame Falls to..." by Sam's Dogcast [box4], "Sakura Rain" by Lily's Garden [box5], and "A Place of Cake" by Zen with Zhang [box6].

Q: I want to listen to Sakura Rain, what should I do?  
A: You can start listening to "Sakura Rain" by tapping on the corresponding picture [box4]

Q: What's the category of the podcasts listed on this screen?  
A: The podcasts listed on this screen are about various topics such as cooking, gardening, and a place of cake. They are categorized under "New & Notable" [box7].

Q: How do I go back?  
A: You can go back by tapping on the "Browse" button [box8]

Q: What's the function of this screen when the user interacts with application?  
A: The screen is for a podcast application where users can browse and play new and notable podcasts, with options to play, download, and search for specific podcasts.