# Phi-Ground Tech Report: Advancing Perception in GUI Grounding

**Miaosen Zhang   Ziqiang Xu   Jialiang Zhu   Qi Dai   Kai Qiu   Yifan Yang**
**Chong Luo   Tianyi Chen   Justin Wagle   Tim Franklin   Baining Guo**

**Microsoft**

## ABSTRACT

With the development of multimodal reasoning models, Computer Use Agents (CUAs), akin to Jarvis from *"Iron Man"*, are becoming a reality. GUI grounding is a core component for CUAs to execute actual actions, similar to mechanical control in robotics, and it directly leads to the success or failure of the system. It determines actions such as clicking and typing, as well as related parameters like the coordinates for clicks. Current end-to-end grounding models still achieve less than 65% accuracy on challenging benchmarks like ScreenSpot-pro and UI-Vision, indicating they are far from being ready for deployment. In this work, we conduct an empirical study on the training of grounding models, examining details from data collection to model training. Ultimately, we developed the **Phi-Ground** model family, which achieves state-of-the-art performance across all five grounding benchmarks for models under $10B$ parameters in agent settings. In the end-to-end model setting, our model still achieves SOTA results with scores of **43.2** on ScreenSpot-pro and **27.2** on UI-Vision. We believe that the various details discussed in this paper, along with our successes and failures, not only clarify the construction of grounding models but also benefit other perception tasks. Project homepage: https://zhangmiaosen2000.github.io/Phi-Ground/

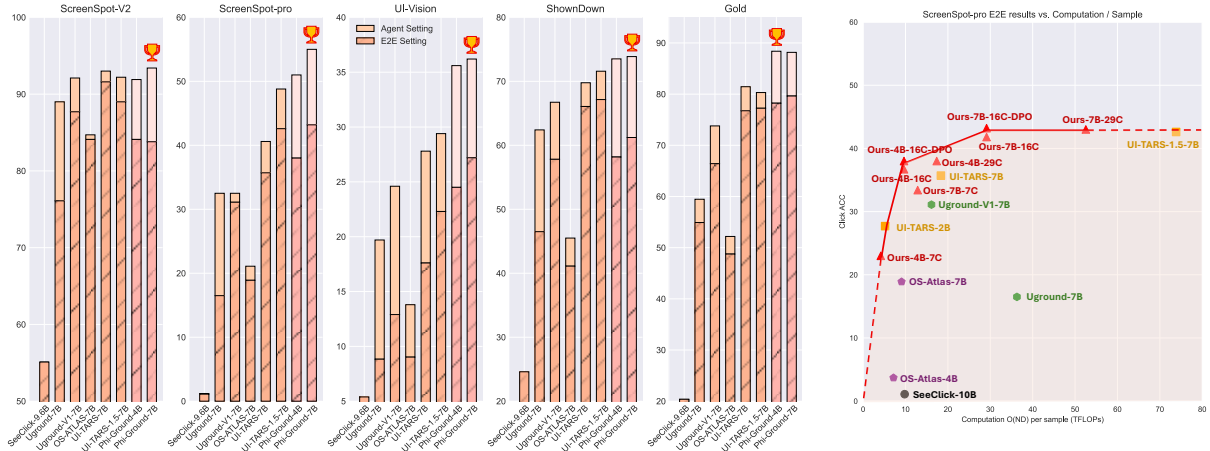**Keywords** GUI grounding · AI agent · Large multi-modal model



Figure 1: **Left**: The comparison chart of our grounding model results across five GUI grounding benchmarks. Our model, trained specifically for the agent setting, achieved SOTA results on all benchmarks under this focus. Even in the general end-to-end model setting, our model attained SOTA results on three of the benchmarks. **Right**: The relationship between model performance and computational cost on ScreenSpot-pro demonstrates that our model supports the Pareto frontier, indicating its efficiency. Most GUI research traditionally considers only the parameter count $N$ for comparison, but our experiments highlight that computational cost during testing, such as the number of image tokens, also significantly impacts performance. The X-axis in the right figure represents $ND$, where $D$ is the number of image tokens. Training and inference latency are more linearly correlated with $ND$ than with $N$. A graph using latency as the X-axis closely resembles the right figure, but latency is often influenced by hardware and acceleration libraries such as vllm, so we did not use latency as X-axis.
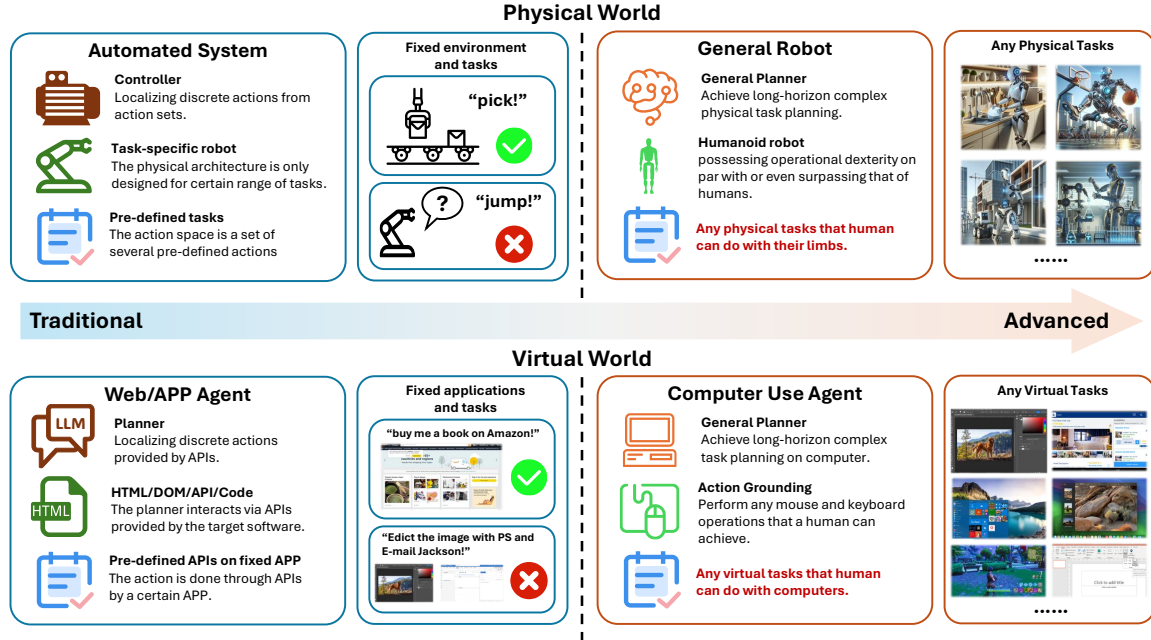
Figure 2: Agent evolution across physical and virtual worlds. Traditional systems rely on fixed controllers and pre-defined workflows to execute domain-specific tasks, either in physical environments (e.g., task-specific robots) or virtual environments (e.g., API-based Web/APP agents). In the modern era, intelligent automation has emerged. In the physical world, general-purpose robots perform versatile limb-based operations. In the virtual world, Computer Use Agents (CUAs) achieve human-level behaviors through general purpose planner, GUI grounding, enabling them to complete any virtual task achievable via mouse and keyboard interactions.

# 1 Introduction

Large model based autonomous agents [1, 2], by leveraging their robust reasoning capabilities and interactions with real-world environments [3], enable humans to tackle more complex tasks and significantly enhance productivity. Given that a substantial portion of modern human work is conducted via computers, Computer Use Agents (CUAs) hold immense potential and commercial value [4, 5]. While CUAs streamline virtual work, robots [6, 7] simplify physical tasks, as shown in Figure 2, and together, they are poised to drive a revolution in productivity. With the development of reasoning models like O3 [8, 9, 10, 11], we are beginning to create prototypes of CUAs, such as OpenAI Operator [12] and Claude Computer Use [13]. However, there remains a considerable gap before CUAs can be fully commercialized. This is due to the irreversible effects of many computer operations [14], such as closing unsaved files, as well as significant concerns regarding user privacy protection [15, 16, 17].

Specifically, CUAs can execute actions through the following methods: APIs (e.g., HTML/DOM, Office scripts), scripting tools (such as AutoHotkey, Power Automate, CLI), and simulated input device interactions (such as mouse clicks and keyboard inputs). Traditional approaches based on APIs [18, 19] and scripts are often constrained by the platform in use, environmental dependencies, and the specific APIs provided by applications [20]. In contrast, solutions based on GUI and simulated input device interactions [21] are aligned with human operations. This alignment not only facilitates user supervision, enhancing privacy and security, but also theoretically allows for any interaction that a human can perform, without being limited by the application. As a result, GUI-based CUAs are becoming a focal point of research in this field, as illustrated in Figure 2.

When accomplishing a task, GUI-based CUA can be divided into two steps: temporal planning and grounding [22, 23]. Planning involves analyzing the task description and the current state to determine the actions that should be taken in the future, while grounding refers to the above mentioned simulation of input devices. Among all interactions, mouse click operations are the most important and common actions for GUI grounding. Since keyboard commands, such as pressing the "A" key, are discrete, MLLMs can effectively handle this type of grounding. Hence, our focus is primarily on mouse commands, where the main challenge lies in the fact that mouse command parameters are screen coordinates, and most MLLMs struggle to accurately identify these coordinates [24, 25, 26]. Therefore, specialized training is required for determining the precise click coordinates.
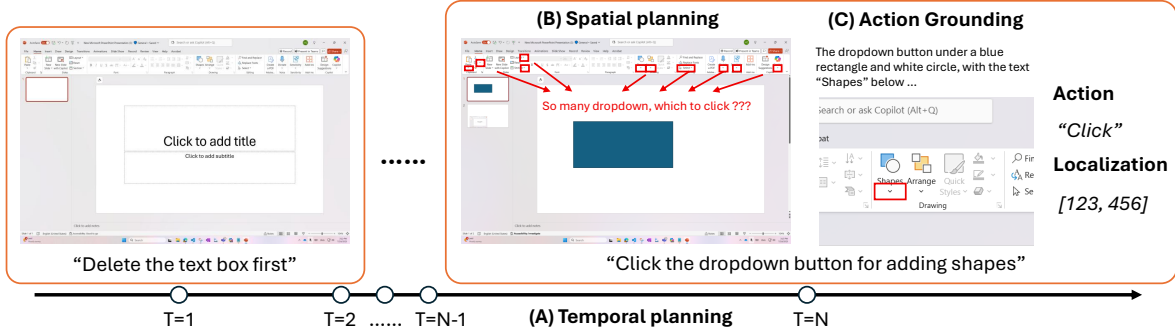
Figure 3: Three levels of task of CUAs. Each coral block represent an action step. We focus on the click action, as it is the most important and common operation. Others such as keyboard input can be effectively handled by MLLM.

Motivated by the above considerations, this paper conducts a detailed empirical study on the training of GUI grounding models. We divide GUI grounding into two components: the first component is spatial planning [27, 23], which involves identifying which specific element in the image needs to be manipulated to execute a given instruction. The second component is localization, where, after determining the target location, the model needs to output the correct coordinates, as illustrated in Figure 3. Many existing CUA models [28, 29, 30] attempt to accomplish the grounding task in an end-to-end manner. However, due to the spatial planning component, which also demands a strong level of expertise, common sense, and spatial reasoning abilities, there is a need for larger or even deep reasoning models. To achieve better results at this stage, we have adopted a two-phase approach [31]. Initially, an advanced MLLM provides a detailed description of the location, followed by our trained grounding model outputting the specific coordinates. In this technical report, we delve into numerous often-overlooked details spanning data, algorithms, and training methodologies. Counterintuitively, we found that many techniques that appear sound and frequently appear in previous work, such as tokenized coordinates [32, 33, 34] , coordinate label smoothing [35, 36], and loss reweighting, become trivial when dealing with large-scale training. We have retained only the following techniques that remain meaningful under extensive training in the main text.

Firstly, at the model input level, we discovered that the order of modality inputs can significantly impact the underlying mechanism of feature modeling, leading to notable differences in results. Secondly, data augmentation [37] is very common in traditional object detection training, yet is rarely mentioned in the era of large models. Our re-experiments indicate that certain data augmentations can greatly enhance results in high-resolution scenarios (such as ScreenSpot-pro). Thirdly, the academic community predominantly focuses on the out-of-distribution generalization ability of models, with insufficient research on how models can appropriately perform in-domain continual learning in specific scenarios (such as considering only Photoshop software). This area, however, holds significant practical value. We discuss data strategies and training algorithms for in-domain post-training. Lastly, existing GUI grounding work typically considers only the size of the parameters when making comparisons. We incorporate the computational load during model inference (primarily influenced by the number of image tokens) into the study of scaling laws[38, 39, 40] and evaluations.

Ultimately, we devised an efficient and rational training recipe, collecting over 40M data samples from multiple sources [21, 31, 41]. By scaling up the training volume, we successfully developed the Phi-Ground model family.

In our evaluation process, to avoid the limitations of relying on a single benchmark and to prevent systemic overfitting [42, 43] due to optimization for a specific benchmark, we consider multiple benchmarks in both ablation studies and final evaluations. We conducted a survey and collected four publicly available test datasets [25, 24, 44, 26], complemented by our internally constructed dataset focusing on commonly used software on Windows, which we refer to as the Gold dataset. This brings our total to five evaluation datasets, whereas most CUA research works [30, 31, 28, 29] on grounding typically utilize only one or two of these. The evaluation results demonstrate that, within the Agent setting we are focused on, our Phi-Ground model achieves state-of-the-art results across all benchmarks, with particularly high scores of 55.0 and 36.2 on ScreenSpot-pro [25] and UI-Vision [24], respectively. Furthermore, in the end-to-end model setting, we also achieved the best results in three of the benchmarks, scoring 43.2 and 27.2 on ScreenSpot-pro and UI-Vision, respectively. These findings indicate that our model possesses strong generalization capabilities.

We also present a detailed erroneous case study in Section 6.2 and the appendix. We consider GUI grounding to be a classic scenario for multimodal model perception [45, 46, 47], and believe that our experiences can be effectively generalized to other fields involving multimodal perception. We hope our research will benefit related domains.
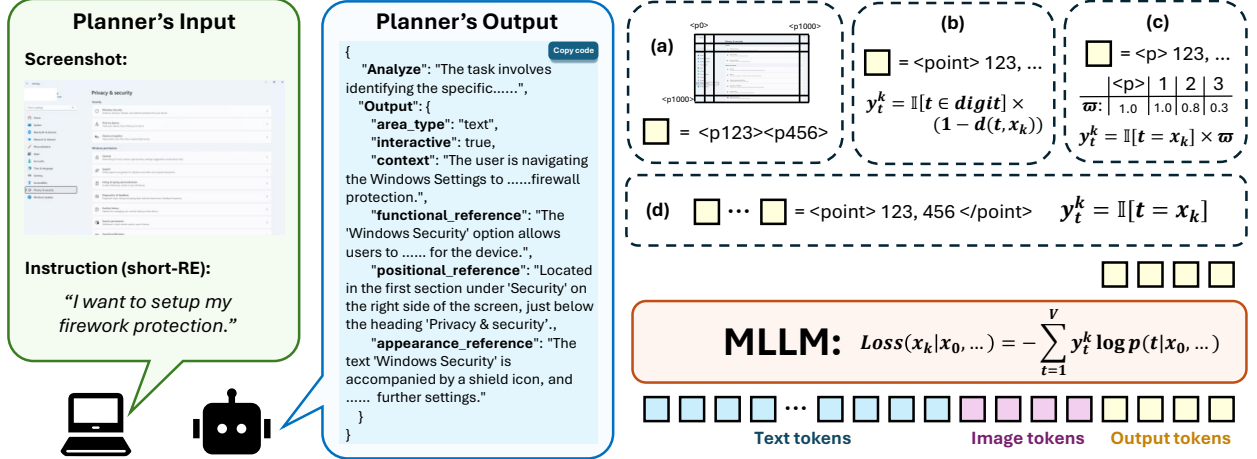
Figure 4: In terms of coordinate representation and loss strategies, we experimented with the following approaches: (a) tokenized coordinate representation, (b) label smoothing for coordinates, (c) loss reweighting, and (d) direct textual representation of coordinates.

## 2 Methodology

At the agent level, we adopt a two-stage implementation approach: a large multimodal model (such as GPT-4O) is utilized to generate detailed and specific Reference Expressions (RE), while a smaller multimodal model that we have trained is responsible for generating coordinates based on the RE. At the model level, we fine-tune a MLLM to directly output coordinates, which are generated in text form. The coordinates are represented as relative values scaled by 1000 and then rounded (e.g., a value of 500 corresponds to a coordinate value of 0.5). We present the following discussion and experiments on the aforementioned implementation.

**Coordinates Representations and Loss**   During the development, we explored several potential implementation approaches as follows:

- **Tokenized coordinates.**   Some existing works represent regions as tokens [32, 33, 34], using learnable special tokens to denote a specific position. However, in GUI grounding, the large image resolution makes using 2D regions as tokens impractical. We attempted to divide both the horizontal and vertical coordinates into 1000 intervals (which is necessary for large screens like 4K displays), with each interval represented by a special token, similar to [48, 49]. Despite experimenting with various initialization and training strategies, our results indicated that introducing a large number of unlearned tokens during finetuning pretrained MLLM can lead to model collapse and poor performance. For further details, see the Appendix C.1.

- **Label smoothing.**   To approximate the classic regression loss, we applied label smoothing to the loss of digit tokens, making the magnitude of the loss proportional to the distance between the prediction and the target. For instance, if the target digit is "5", we assign smaller labels to "4" and "6" to indicate their proximity to the target. The detailed implementation can be found in the Appendix C.2. However, the experimental results indicate that this technique may offer improvements only when the training dataset is small and the batch size is minimal. When we eventually increased the batch size to 2048 and the total training samples reached million level, the application of this technique showed no significant impact or improvement.

- **Loss re-weighting.**   Following a similar idea as Label smoothing, we assigned different loss weights to the digits representing the hundreds, tens, and units places of the coordinates. The results were consistent with the previous findings: not only were the corresponding parameters extremely sensitive, but the technique also lost its advantages as the batch size and training volume increased. See Appendix. C.3.

From the above experiments, we found that the most straightforward text output combined with GPT loss (NTP) can effectively facilitate training. We chose to use relative coordinates for expressing positions because their value range is fixed and relatively dense (as the resolution width and height of GUIs are generally greater than 1000 pixels), which facilitates better model learning. Additionally, since relative coordinates are floating-point numbers that are $\leq 1$, the leading "0." is insignificant. Therefore, following [31], we scaled these values by multiplying them by 1000.
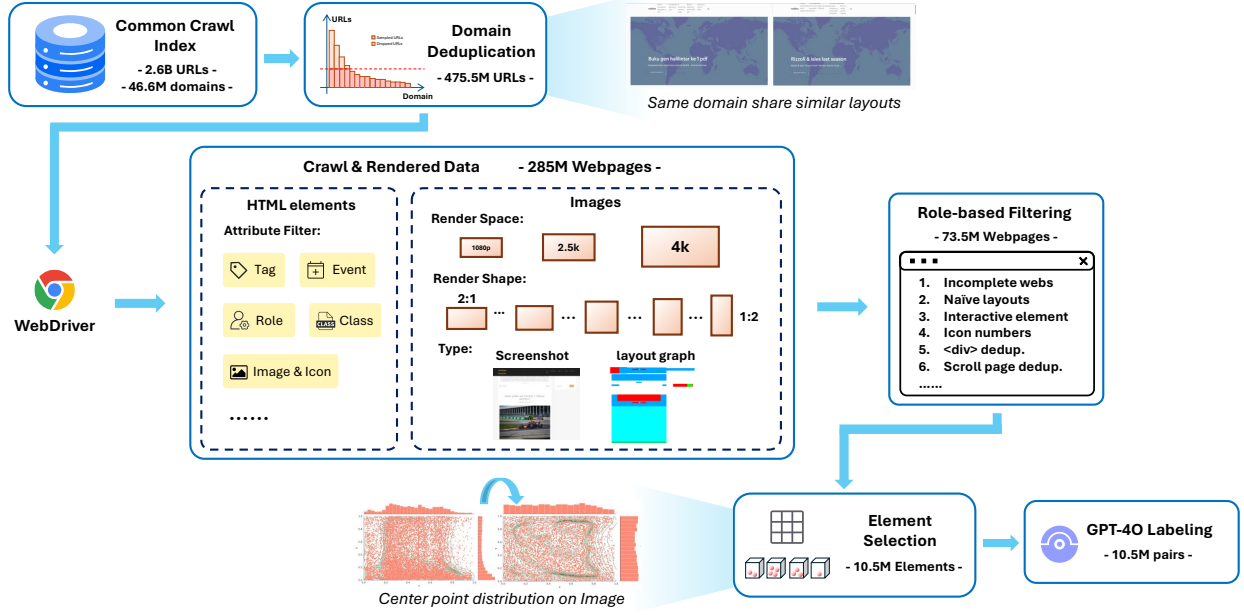
Figure 5: CommonCrawl data processing pipeline.

# 3 Data Preparation

## 3.1 Processing Open Source Data

We utilized a portion of the open-source data from OS-Atlas [31], which includes several subsets such as Windows, macOS, Linux, and Fineweb. Additionally, we incorporated data from SeeClick [26], E2ISynth [41] and GUIAct [21]. The total data volume amounts to approximately 10M entries. By employing the prompt outlined in Appendix E.1, we re-annotated all the data using GPT-4O to generate "Long-gold" type reference expressions (as described in Sec 4). Furthermore, we constructed training data by randomly combining three types of reference expressions.

## 3.2 CommonCrawl Data

To acquire larger-scale data for better scaling up of training, we also obtained web pages from CommonCrawl [50] and rendered screenshots to generate training data. However, the web data contained a significant amount of noisy data that caused training failures. To address this, we constructed a highly specific data cleaning pipeline, as illustrated in Figure 5. Below are the detailed steps of each stage:

**Index and domain deduplication** We utilized the *CC-MAIN-2024-46* crawl from CommonCrawl. After a basic deduplication of URLs (exact match), filtering by language (retaining only English), and webpage status (retaining only 2xx, 301, and 302), we were left with 2.6 billion URLs. These 2.6B URLs originate from 45.6 million unique domains, with the number of pages from the same domain displaying a long-tail distribution. For instance, the largest domain contains 204K different pages. We observed that pages from the same domain exhibit strong consistency in layout. Therefore, to ensure the generalizability of our model, we performed random sampling so that no more than 50 pages were selected from each domain. After this round of sampling, we were left with 475.45M URLs.

**Rendering** We utilized the Selenium library and Google Chrome Driver to render webpage screenshots. During the rendering process, we randomly selected from three different pixel areas corresponding to 1080p, 2K, and 4K screen resolutions. The aspect ratio of the images was randomly chosen between 2:1 and 1:2. For the elements within the webpage HTML, we designed several rules for filtering and retaining them, as detailed in Appendix D.1.1. This process allowed us to preserve elements that are likely to be interactive components. At this stage, we save webpage screenshots, element information, and layout graphs (with different types such as interactive text buttons, interactive icon buttons, and images corresponding to specific colors). After this stage, there retained 285M webpages.

**Rule-based filtering**     Subsequently, we designed more fine-grained filters and deduplication techniques at the webpage and element levels based on the preserved webpages. The specific details are provided in Appendix D.1.2. These filters eliminated many erroneous and overly simplistic webpages. After this phase, 73.5M webpages remained.

**Element selection and labeling**     Finally, when selecting elements, we consider the distribution of element centroids and their types, as we found it necessary to do so in Sec. 5.3. Specifically, we discretize and uniformly sample across various regions of the canvas (see Appendix D.2). During sampling in a discrete area, we prioritize sampling icon elements, as they are less frequent. We sample only one element per webpage screenshot. Consequently, after this stage, 10.5M elements and screenshots remain. Finally, we use GPT-4O to annotate all the data with prompt in Appendix E.1.
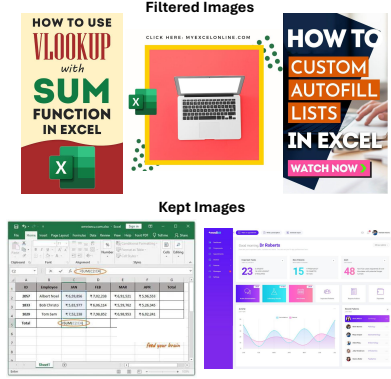


Figure 6: Filtered and kept image of the Web Search Data with the classifier.

Table 1: The applications distribution of the search queries. These applications are derived from an internal statistical list of the most frequently used applications on Windows.

| Domain | # Apps | Example apps | # Queries |
|---|---|---|---|
| Office & Productivity | 10 | Word, Excel | 50 |
| Media Playback & Creation | 11 | VLC, Photoshop, Clipchamp | 55 |
| Gaming Platforms | 10 | Steam, League of Legends | 50 |
| Social & Communication | 10 | Discord, Zoom | 50 |
| Security & System Tools | 10 | CCleaner, Malwarebytes | 50 |
| Core Windows UI | 15 | Settings, File Explorer | 75 |
| **Total** | **66** | - | **330** |

## 3.3   Web Search Data

We construct a complementary, high-resolution screenshot corpus with the Bing Image Search v7 API. We first generate queries from 6 UI domains which covering a total of 66 desktop applications, as shown in Table 1. For each application we manually construct file search phrases such as *"Word ribbon interface"* or *"screen shot of VLC playback controls"*, producing 330 queries altogether. Every query retrieves 2048 candidate images and enforces a minimum resolution of $200 \times 200$ px. **The request also specifies Bing's *license* flag to ensure that all the images we use comply with copyright requirements**. Each downloaded image is then scored by a CLIP-based classifier to remove non-screenshot images from the search result. For the final filtered screenshots, we employed OmniParserV2 [51] to annotate all bounding boxes. We then used the same method described in 3.2 to label reference expressions and perform rule-based filtering at the bounding box level. This process resulted in approximately 158K data samples.

## 3.4   Human Labeled Data

To address our focus on specific scenarios (Windows and common applications) and to explore in-domain training techniques, we have developed a pipeline for constructing human-labeled data. The data construction process involves three steps:

- In the first step, labelers use custom screen recording software to interact with the target scenarios, which include specific software or Windows system settings pages. They are required to navigate and operate various pages and functions of the software. The screen recording software automatically captures different pages and retrieves the UIATree of the page. If the software does not implement UIA, we utilize OmniParser-V2 [51] to obtain the bounding boxes of elements.

- In the second step, recognizing that the acquired bounding boxes contain numerous errors, redundancies, and non-interactive content (such as text portions in Word documents), we developed another software tool to enable labelers to remove erroneous bounding boxes.

- Finally, we employ the Long-gold method, as described in Sec. 4.2, to have GPT-4O annotate references for all bounding boxes.
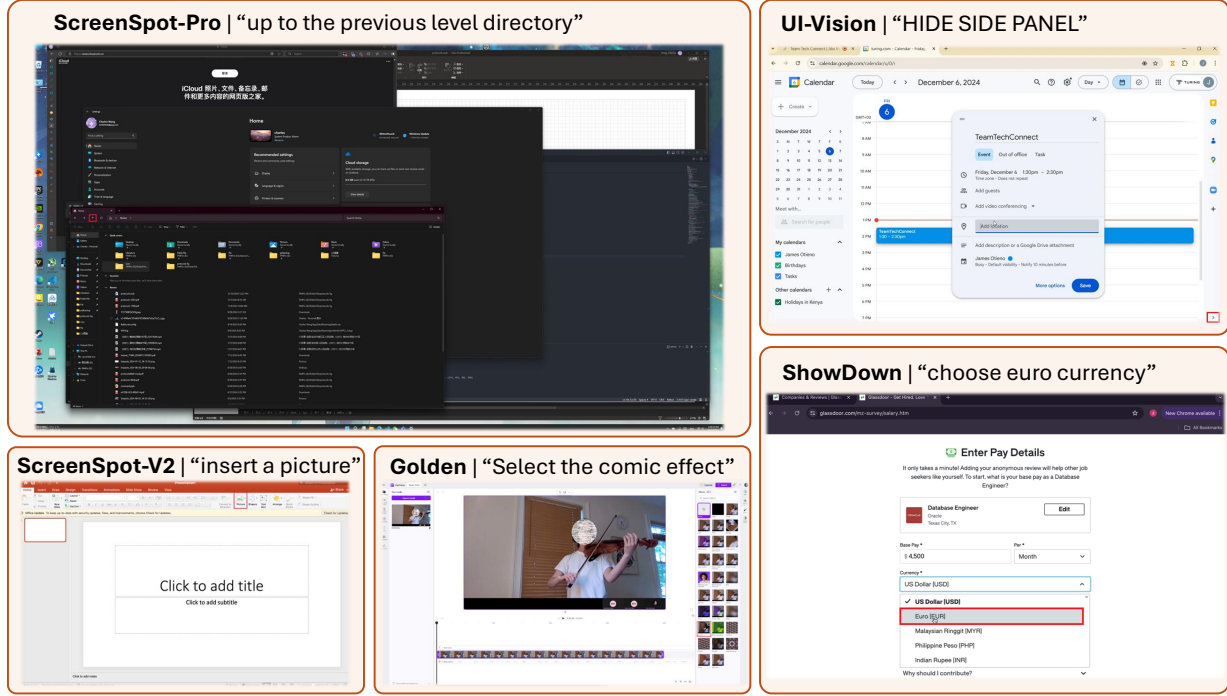
Figure 7: Examples of benchmarks used in evaluation.

Ultimately, we generated 80K training samples derived from a comprehensive suite including Microsoft Office, Windows settings, and over a dozen common software applications such as 7zip, PhotoShop, ClipChamp and audio tools.

## 4 Evaluation Settings

### 4.1 Benchmarks

To ensure the model's generalization capability and to avoid systematic overfitting to well-known benchmarks such as ScreenSpot, we have gathered several recent open-source and internally developed evaluation datasets. This approach aims to ensure the comprehensiveness of our testing.

**ScreenSpot V1 & V2**   ScreenSpot [26] is the first realistic GUI grounding benchmark that includes over 600 interface screenshots from mobile, desktop, and web environments. It encompasses two different types of elements, namely, texts and icons. However, due to several incorrect or ambiguous annotations in the original benchmark, OS-Atlas [31] introduces an updated version called ScreenSpot-V2, which corrects various mislabeled instances in the original ScreenSpot dataset.

**ScreenSpot Pro**   ScreenSpot-Pro [25] focuses on high-resolution GUI tasks in professional applications. It contains 1,581 samples encompassing 23 applications across 3 operating systems, and the micro-elements only cover average 0.07% of the screen, which leads to a highly challenging assessment.

**UI-Vision**   UI-Vision [24] provides three fine-to-coarse grained tasks—Element Grounding, Layout Grounding, and Action Prediction. Similar to ScreenSpot-pro, this benchmark encompasses another set of practical applications such as education and entertainment. It also categorizes three different types of references for testing purposes. Due to the selection of smaller buttons or those requiring more expert knowledge, this benchmark is equally challenging as ScreenSpot-pro. Moreover, it features a larger sample size, enhancing the stability and reliability of the evaluations.

**Showdown-click-dev**   Showdown [44] is a collection of 5,679 left clicks of humans performing various tasks in a macOS desktop environment. A subset containing 557 clicks was released on GitHub and huggingface, namely

Showdown-click-dev. It reported grounding results and latency of several advanced CUAs like OpenAI Operator [12] and Claude computer use [13].

**Gold Dataset (proprietary)**   We have also internally constructed an evaluation dataset tailored for application scenarios within Windows. This dataset includes six scenarios: Photoshop (213 samples), ClipChamp (179 samples), PowerPoint (82 samples), Excel (107 samples), Word (87 samples), and Windows settings (1266 samples). The bounding boxes and instructions (Short RE) for each scenario were annotated by professional engineers familiar with Windows. While these share the same domain as the human-labeled data in Sec 3.4, the labelers come from different teams and use different computer settings, such as theme colors and screen sizes. In contrast, the selection of elements and the generation of REs in our training set's human-labeled data are automated or AI-generated, lacking expert knowledge. We employ this Gold dataset to evaluate in-domain training research. Additionally, it serves as an effective validation of practicality within Windows scenarios. In subsequent experiments, to facilitate faster testing, we manually extracted a smaller test set of 211 samples from this dataset, referred to as **Gold-S**.

We present examples of these benchmarks in Figure 7. In addition, there are benchmarks like WinClick [52] that are still in the process of being open-sourced. As a result, we are unable to present their results in this paper. However, we plan to showcase these results on the project homepage once they are released.

## 4.2   Evaluation Protocols and metrics

**Reference expression type**   A Reference Expression (RE) refers to a segment of referring text provided to a grounding model or GUI agent model concerning a specific interactive area, such as a button. This expression can be an indirect instruction like "close the webpage", or a direct and specific description such as "the blue settings icon in the upper left corner". The use of different REs during training and testing can significantly impact the model's performance. Generally, we desire the model to generalize well across various REs. However, when using indirect REs, the model also needs some planning capability. To decouple model capabilities of planning and perception for research and enable efficient application in smaller models, previous works [31] introduced a simplified agent setting: employing more powerful MLLMs like GPT-4O [53] and O4-mini [8] for planning and generating more detailed REs, while a smaller grounding model is used to produce coordinates. In this paper, to minimize misunderstanding, we define the following three types of REs.

- **Short / instruction.** "Short RE" refers to a more concise form, potentially containing some instructions that require planning. During testing, we use "Short RE" to denote the reference provided by the benchmarks. This type of RE will be represented as the ***"End-to-end model setting"*** in the subsequent results tables.

- **Long / agent.** Following [30, 31], for short RE, we can utilize advanced MLLM to expand it into three more explicit and detailed types of RE: functional, positional, and appearance. We allow the large model to generate these three types of references and concatenate the references together, forming what we call a Long RE. During testing, the input to the MLLM consists of only the screenshot and the short RE, and the output is the Long RE, which is then passed to a smaller model to generate coordinates. In subsequent tables, this approach is referred to as the ***"Agent setting"***. The specific system prompt for generating Long RE can be found in Appendix E.2.

- **Long-gold.** In the process of constructing the training dataset, on one hand, there is no short RE provided, while on the other, we can utilize ground truth (GT) bounding boxes. We also instruct the MLLM to generate the aforementioned three types of long REs. The difference lies in the input we provide to the MLLM, which includes a screenshot annotated with the GT bbox and a cropped image of the target region. This allows the model to better observe and generate high-quality REs. However, since the generation process relies on the GT, this Long-gold RE is used solely for generating training data and not for evaluation. The specific system prompt can be found in Appendix E.1.

Since this paper primarily focuses on model perception rather than planning, and all training utilizes Long-gold RE, unless otherwise specified, all experimental results in Section 5 are based on Long RE generated by GPT-4O.

**Metrics**   In all tables within this paper, unless otherwise specified, the reported metric is click accuracy. The benchmark provides a GT bounding box, and the model being tested generates a click coordinate. If the model's output format is a bounding box, the center of the box is taken as the click coordinate. A click is considered correct if it falls within the area of the GT bounding box; otherwise, it is deemed incorrect. This method is used to calculate the accuracy.
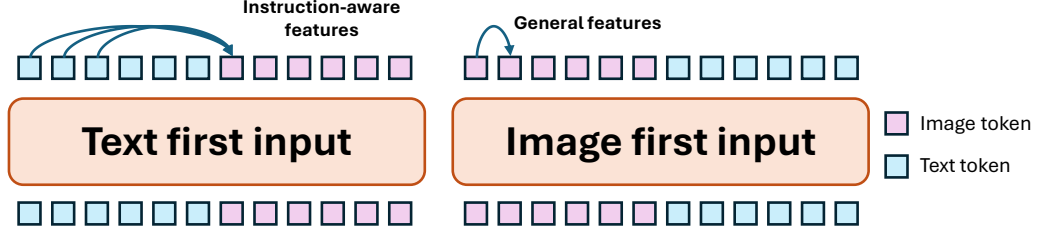
Figure 8: Illustration of the impact of modal input order on model training.

# 5 Experiments

## 5.1 Input / Output Format

**The order of input modality.** We investigated and experimented with the order in which text (or RE) and images are input into the model. The results displayed in Table 2 indicate that inputting text before images yields significantly better outcomes than the reverse order, aligning well with our expectations. Figure 8 illustrates the underlying mechanism: due to the use of causal masks in transformer decoder's attention, earlier tokens cannot be updated using later tokens. When images are input into the model first, the image tokens in the later layers, adapted through training, can be understood as being tailored for grounding tasks but unrelated to the RE. Conversely, when the order is reversed, the modeling of image tokens becomes instruction-aware. For perceptual tasks, the effectiveness of image modeling directly influences the final results. This observation bears significant resemblance to existing research related to instruction-aware models [54, 55, 56]. A simple modification of the modality order can effectively achieve this outcome.

Table 2: Comparison of input order of modalities.

| Input format | max_crops | SSv2-Mobile | SSv2-Desktop | SSv2-Web | SS-pro | Gold-S |
|---|---|---|---|---|---|---|
| Image first | 6 | 85.1 | 82.7 | 81.1 | 14.6 | 84.3 |
|  | 15 | 85.4 | 83.2 | 82.6 | 27.7 | 86.1 |
| Text first | 6 | 87.5 (+2.4) | 83.2 (+0.5) | 83.1 (+2.0) | 18.7 (+4.1) | 86.7 (+2.4) |
|  | 15 | 87.1 (+1.7) | 85.3 (+2.1) | 84.2 (+1.6) | 30.6 (+2.9) | 89.2 (+3.1) |

**Output format.** We investigated the impact of output formats on results, where output formats refer to the representation of coordinates, such as bounding boxes and points. Depending on the scenario, developers might need to employ different output formats. For instance, in scenarios involving the implementation of an OS agent, grounding only requires providing the coordinates for a click, which means delivering a point. However, some developers may wish the model to generate a bounding box. For the output format of boxes, using the center point of a box can serve to determine the click position, thereby accommodating the previous scenario. We discussed the following formats and implementation strategies:

- **Point+OmniParser.** The model will output text in the following format: $< point > mid_x, mid_y < /point >$, where $mid_x$ and $mid_y$ represent the coordinates of the central point. In scenarios involving box output, we first utilize OmniParser-V2 [51] to annotate all elements within the image. Subsequently, we select the target box using the points provided by the grounding model.
- **XYXY.** The model will output a bounding box with format: $< box > x_1, y_1, x_2, y_2 < /box >$, where the $(x_1, y_1)$ and the $(x_2, y_2)$ are the top-left and bottom-right point.
- **XYWH.** The model will output a bounding box with format: $< box > x_1, y_1, w, h < /box >$, where the $(x_1, y_1)$ is the top-left corner and $w, h$ are the width and height of the box.
- **MidWH.** The model will output a bounding box with format: $< box > mix_x, mid_y, w, h < /box >$, where the $(mid_x, mid_y)$ is the center point of the bounding box and $w, h$ are the width and height of the box.

Results are shown in the table 3. We found that directly outputting the point format achieved the best Click ACC, while the MidWH format better balanced Click ACC and the precision of the detection box. These conclusions align well with intuition. In terms of output precision for detection boxes, the XYXY format demonstrated the best performance. We observed that integrating OmniParser-V2 did not yield better results. On the one hand, the performance on the

Table 3: Results of training using different output formats. Note that for the point+OP format, the Click ACC reported in the table is derived directly from using the model's output point as the click target. In contrast, IoU-related metrics are calculated by selecting a box annotated by OmniParserV2 based on the model's output. The results in parentheses represent the average of the SSv2-desktop and SSv2-web subsets.

| Output format | ScreenSpot-V2 | | | | | ScreenSpot-pro | | |
|---|---|---|---|---|---|---|---|---|
| | Click ACC | IoU@0.3 | IoU@0.5 | IoU@0.8 | IoU | Click ACC | IoU@0.5 | IoU |
| point+OP | **85.5** | 47.3 (61.8) | 39.3 (53.8) | 22.0 (**30.9**) | 36.7 (48.1) | **30.6** | 22.3 | 18.8 |
| XYXY | 84.7 | 72.4 (71.2) | **59.6 (59.0)** | 26.6 (28.3) | **53.7 (53.7)** | 30.0 | **23.8** | **20.1** |
| XYWH | 84.2 | 73.6 (70.9) | 58.9 (57.7) | 22.5 (25.2) | 52.6 (52.1) | 29.4 | 21.8 | 18.6 |
| MidWH | 85.0 | **73.6 (71.5)** | 57.1 (58.1) | 21.7 (25.0) | 51.6 (52.2) | 30.4 | 22.4 | 19.2 |

ScreenSpot-mobile subset was exceptionally poor. On other subsets, the results were still slightly lower compared to models that output box coordinates end-to-end. However, for high box precision required situations (IoU > 0.8) and high-resolution scenarios like ScreenSpot-Pro, the gap was reduced, indicating that OmniParser's implementation may have certain advantages in high-resolution contexts and generate more accurate boxes.

## 5.2 Data Augmentation

---

**Algorithm 1** random crop

---

**Require:** $img, box$,Croping probability: random_crop, Minimal cropping ratio: min_crop
**Ensure:** Cropped image $img$,Updated box $box$
1: **if** rand() < random_crop **then**
2: $\quad (w, h) \leftarrow img.size$
3: $\quad x_{left} \leftarrow box[0]$
4: $\quad x_{right} \leftarrow 1.0 - box[2]$
5: $\quad x_{crop\_factor} \leftarrow$ min_crop + rand() $\times$ (1 − min_crop)
6: $\quad crop_{x1} \leftarrow w \times x_{left} \times (1 - x_{crop\_factor})$
7: $\quad crop_{x2} \leftarrow w \times (box[2] + x_{right} \times x_{crop\_factor})$
8:
9: $\quad$ ... Same for Y axis ...
10:
11: $\quad img \leftarrow$ CropImage($img$,
12: $\quad\quad (crop_{x1}, crop_{y1}, crop_{x2}, crop_{y2}))$
13: $\quad box \leftarrow$ UpdateBox($box$)
14: **end if**

---

**Algorithm 2** random resize and padding

---

**Require:** $img$, Target dimensions $(W, H)$, Random resize probability: random_resize, maximin screen size: max_screen_size
**Ensure:** Output image $canvas$, Updated box $box$
1: $(w, h) \leftarrow img.size$
2: $canvas \leftarrow$ WhiteImage($W, H$)
3: **if** rand() < random_resize **then**
4: $\quad S_{max} \leftarrow \min(1.0, \frac{W}{w}, \frac{H}{h})$
5: $\quad S_{min} \leftarrow \frac{W}{\text{max\_screen\_size}} \times S_{max}$
6: $\quad scale \leftarrow S_{min} + \text{rand()} \times (S_{max} - S_{min})$
7: $\quad (\hat{w}, \hat{h}) \leftarrow (w \times scale, h \times scale)$
8: $\quad img \leftarrow$ Resize($img, (\hat{w}, \hat{h})$)
9: $\quad pos \leftarrow (\text{randint}(0, W - \hat{w}), \text{randint}(0, H - \hat{h}))$
10: **else**
11: $\quad$ Scale the image to fit in canvas...
12: $\quad pos \leftarrow (0, 0)$
13: **end if**
14: $canvas$.paste($img, pos$)
15: $box \leftarrow$ UpdateBox($box$)

---

We investigated the effects of two data augmentation techniques on the training of UI grounding models:

**Random Crop** When a user's software interface is larger than their screen, or under certain interface scaling conditions, it may result in the display of an incomplete page. Additionally, the ScreenSpot contains a small number of incomplete UI. Motivated by these observations, we introduce random crop as a data augmentation technique. We define the right edge of the bounding box and the right edge of the image as coordinates 0 and 1, respectively. We then select a random number within the range [min_crop, 1], which determines the new right boundary of the image. The same process is applied to other directions. Notably, we use the same random number for both the left and right sides (as well as the top and bottom), ensuring that the cropping is proportional. This approach maintains the positional integrity of objects, such that if a box is located on the left side of the original image, it will remain on the left side of the cropped image. This helps avoid potential errors or changes in positional references. The detailed implementation is shown in Algorithm 1.

**Random Resize**   When a user reduces the software interface size or when the screen resolution is very high, items may become very small. To address this issue, we introduce a random resize data augmentation technique. The core of this data augmentation method is to shrink the image and place it onto a fixed-size white canvas. The canvas size is generally related to the design of the LMM model, particularly how it partitions the image into patches. We set a maximum screen size (e.g., 4K), and the canvas size is randomly selected between the size of the images in the training set and the maximum screen size. The purpose of this approach is to leverage the inherent size of the training set images (for instance, if the training set images are already large, excessive resizing should be avoided). For more details, please refer to Algorithm 2.

Table 4: Data augmentation ablation experiments.

| Data Augmentation | SSv2-desktop | SSv2-web | SSv2-mobile | SS-pro | Gold-S |
|---|---|---|---|---|---|
| base | 85.9 | 84.7 | 88.1 | 24.8 | 90.0 |
| base + R-crop | 86.2 | 84.9 | 87.9 | 23.6 | 89.8 |
| base + R-resize | 85.9 | **85.1** | **88.3** | **32.8** (+8.0) | 89.1 |
| base + R-crop + R-resize | **86.5** | **85.1** | **88.3** | 32.7 (+7.9) | **91.0** |

We conducted several sets of hyperparameter experiments, and the optimal combination is as follows: The probability for Random Resize is set to 100%, with a maximum screen size of 4096, while the probability for random cropping is 0.3 with a min_crop = 0.7. We utilized these hyperparameters in all subsequent experiments. Table 5 presents the results of the ablation study on data augmentation. The results indicate that in high-resolution testing environments such as ScreenSpot-pro, employing random resize significantly improves performance. On the other hand, random cropping does not have a substantial impact in various scenarios.

## 5.3   Data Distribution



Figure 9: Data distribution.

Table 5: Data ablation experiments. "Bing" refers to our BingSearch dataset, "CC" denotes the CommonCrawl dataset, and "CC" and "CC-re" represent the datasets before and after applying our proposed re-sampling technique, respectively.

| ID | Data | SSv2 | SS-pro | Gold-S |
|---|---|---|---|---|
| 1 | base | 85.3 | 32.3 | 89.1 |
| 2 | 97%base + 3%Bing | 85.5 | 32.7 | 89.6 |
| 3 | 70%base + 30%CC | 85.8 | 31.5 | **90.0** |
| 4 | 70%base + 30%CC-Re | 85.5 | 33.1 | 89.6 |
| 5 | 67%base + 3%Bing + 30%CC-Re | **86.1** | **33.3** | 89.7 |

We analyzed the distribution of the relative positions of the center points of all detection boxes within each dataset, as illustrated in Figure 10. Our findings reveal that data from different sources exhibit distinct distribution patterns. For instance, we observed that web-rendering data such as FineWeb, SeeClick, and our CommonCrawl data (when not resampled) exhibit an identical distribution pattern, as shown in the first graph of Figure 10(a). These datasets almost completely lack buttons on the right side, which can be attributed to the common web design practice of placing sidebars on the left and reserving the right side for scroll bars. However, such a distribution is not universal. For instance, in other datasets, the distribution tends to be more uniform. Our BingSearch demonstrates the best uniformity and diversity. Different software exhibits its unique distribution, as evidenced by our human-labeled data.

Therefore, we proposed an algorithm to resample our own CommonCrawl data. The basic idea of the algorithm is to divide the image into a $50 \times 50$ grid and sample a fixed number of points from each grid cell. This approach ensures that the central points are uniformly distributed by area, as detailed in Appendix D.2. During the resampling process, we introduced a sampling factor to balance the trade-off between sampling rate and uniformity, as illustrated in Figure 10(b).

To validate the generated data and the effects of the aforementioned data resampling, we first constructed a base setting using open-source data, as shown in Figure 9. The data proportions are roughly proportional to the size of the subset.

**(a) The center point distribution of different training datasets**



**(b) The center point distribution and remaining data proportion (RES) when changing sampling factor (R)**



**(c) The center point distribution of different benchmarks**



Figure 10: Center point distribution of training and evaluation data.

We sampled 5M data points from this distribution to train the model. Based on this, we gradually replaced with our data, such as replacing 3% of the BingSearch data (3% of 5M is approximately equivalent to 1 epoch of BingSearch) or 30% of the CommonCrawl data.

The results are shown in Table 5. Comparing experiments ID-1, ID-2, and ID-3, the model's performance on specific benchmarks improved after incorporating our data. Since we primarily replaced web data (Fineweb and SeeClick), this indicates that our data surpasses the web data from open-source datasets, mainly due to the greater variety of our data. Comparing experiments ID-3 and ID-4, there is a significant improvement on the ScreenSpot-pro dataset. This suggests that a uniform distribution of boxes is more beneficial for generalization in high-resolution scenarios. Finally, combining all our data achieved the best results (ID-5).

## 5.4 In-domain Post-training

In this section, we will utilize our human-labeled data and explore fine-tuning strategies for target scenarios. In practical applications, developers might have a small set of target software that they wish to cover with their agent. We will use Adobe Photoshop (PS) as an example. Our human-labeled data includes screen recordings processed using the method described in Sec. 3.4. The bounding boxes are derived from UIA-tree, while the references are sourced from GPT-4o. For our evaluation set, we focus on the PS subset of our Gold dataset, as well as the PS subset from ScreenSpot-pro. It is important to note that our Gold dataset was provided by another team, and both the bounding box and reference annotations were done by humans. Therefore, there may be discrepancies with our training set.

**Domain Finetuning** When performing domain fine-tuning, we consider the following strategies:

- Strategy A involves directly incorporating domain-specific data into the pre-trained model.
- Strategy B entails performing SFT of a pre-trained model (without using domain data) using the domain data.
- Strategy C first introduces a small proportion of domain data during the pre-training phase, followed by using a larger proportion of domain data during the SFT phase.

We set the pre-training data size to 5M and the SFT data size to 200K. Training a subset for too many epochs (e.g., more than 10) can lead to overfitting and other issues. Thus, we limit the number of epochs of domain data for both the pre-training and SFT stages to approximately 3 epochs. Due to the difference in total data volume, the proportion of domain-specific data will vary accordingly.

Table 6: Comparison of data selection strategies during in-domain post-training and pretraining.

| Strategy | Pretrain | | Finetune | | General | | Domain | |
|---|---|---|---|---|---|---|---|---|
| | base | PS (r / ep) | base | PS (r / ep) | SSv2 | SSpro | Gold-PS | SSpro-PS |
| - | 100% | 0% / 0 | - | - | 86.1 | 33.3 | 68.1 | 39.2 |
| A-1 | 99.5% | 0.5% / 2.2 | - | - | **86.2** | 32.7 | 71.8 | 41.2 |
| A-2 | 99.0% | 1.0% / 4.5 | - | - | 86.0 | **34.4** | 72.7 | 41.2 |
| B-1 | 100% | 0% / 0 | 88.0% | 12.0% / 2.1 | 84.2 | 31.4 | 72.0 | 41.2 |
| B-2 | 100% | 0% / 0 | 76.0% | 24.0% / 4.3 | 84.4 | 30.7 | 74.2 | 41.2 |
| C | 99.5% | 0.5% / 2.2 | 88.0% | 12.0% / 2.1 | _85.8_ | _34.2_ | **74.8** | **43.1** |

The results, as shown in Table 6, indicate that directly incorporating domain-specific data into the training set has a minimal impact on the model's general capabilities. However, the improvement in performance for the specific domain is quite limited. This may be due to the relatively small proportion of domain-specific data, which constitutes only a minor part of the optimization objective. In the case of Strategy B, where domain-specific data is introduced only during the fine-tuning phase, we observed a significant improvement in in-domain results. However, we also noted a marked decline in the model's general capabilities after fine-tuning. This suggests that incorporating previously unseen domain data during fine-tuning may lead to overfitting and catastrophic forgetting of pre-trained knowledge.

Strategy C effectively balances general capabilities and in-domain performance. It maintains strong general abilities from the pre-training phase while also achieving the best results within the domain. During our fine-tuning process, we found that the gradient norm for Strategy C was smaller than that for Strategy B. This suggests that the decline in general capabilities observed in Strategy B may also stem from training instability caused by the introduction of new data. Since Strategy C is exposed to domain data during the pre-training phase, it experiences a more tempered impact during fine-tuning.

**Post-training Algorithms** Beyond the impact at the data level, we also explored the effects of training algorithms during the post-training phase. We primarily considered the following categories of algorithms: Supervised Fine-Tuning, Curriculum Learning, and Reinforcement Learning. For reinforcement learning, we further examined several algorithms, including reject sampling finetuning (which is not strictly traditional reinforcement learning but has been shown by many studies [57, 58, 59, 60] to have similar characteristics and effects), DPO [61], REINFORCE [62], PPO [63], REINFORCE++ [64] and GRPO [65]. However, we failed to make positive increase on results with PPO, REINFORCE++, and GRPO, primarily due to three reasons: (1) The pure perception tasks lacked textual exploration and reasoning and hence, in this situation, the essence of using RL loss will be more closely aligned with the original SFT loss. (2) Although some previous work [66, 67, 68] has shown that RL can provide benefits in purely perceptual tasks, these studies typically begin with relatively low baselines or apply RL directly from scratch without prior SFT. In contrast, our work involves pre-training on in-domain tasks to an optimal level, leaving little room for further optimization. (3) The absence of exploration and reasoning led to low diversity in the answers among rollouts for the same sample, frequently resulting in all rollouts being either entirely correct or entirely incorrect (about 70% probability). This made the training of algorithms like GRPO trivial and further hindered the effective training of the critic model in PPO. This prompted us to consider using simpler RL algorithms such as reject sampling finetuning and DPO, as well as manual intervention in the selection of rollout samples.

We employed a unified framework, as illustrated in Figure 11, to integrate SFT, DPO, curriculum learning, and reject sampling finetuning. We ensured that the total number of samples used was consistent with the 200K samples used in SFT. Taking DPO as an example, the effective training samples only include "Case 1" in Figure 11, resulting in less than 100K samples being used for training. We implemented DPO using the trl library [69], which also includes several variants of DPO, such as sigmoid [61](original DPO implementation), hinge [70], IPO [71], exo [72], nca [73],
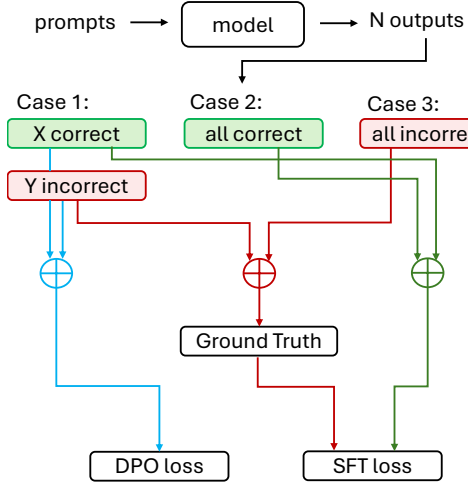
Figure 11: ■:DPO. ■:Curriculum learning. ■:Reject sampling finetuning. $X + Y = N$

Table 7: Results of different post-training algorithms which further optimized a pre-trained model that had reached saturation. To our surprise, we found that for purely perceptual tasks, DPO could still enhance the results of the pre-trained model. However, this conclusion primarily stems from data-level adaptability and robustness, which is fundamentally different from the mechanisms of LLM reasoning.

| Algrithm | General | | | Domain | |
|---|---|---|---|---|---|
| | SSv2 | SSpro | Gold-S | Gold-PS | SSpro-PS |
| Base Model | 86.2 | 32.7 | 88.5 | 71.8 | 41.2 |
| SFT | 85.8 | 34.2 | 89.1 | 74.8 | 43.1 |
| Curriculum-L | 86.1 | 33.7 | 89.4 | 75.1 | 45.1 |
| Reject Sampling | *86.2* | 34.5 | 88.6 | 75.5 | 45.1 |
| DPO-sigmoid | *86.2* | **35.2** | *90.5* | **76.8** | **49.0** |
| DPO-IPO | **86.4** | *34.8* | 90.3 | 75.5 | 43.1 |
| DPO-NCA | 85.2 | 33.9 | 88.7 | 75.7 | 47.1 |
| DPO-DiscoPoP | 85.4 | *34.8* | **91.0** | 76.3 | **49.0** |

robust-dpo [74], SPPO [75], AOT [76], DiscoPoP [77], and APO [78]. For each variant of DPO, as well as for SFT, curriculum learning, and reject sampling, we conducted hyperparameter tuning. This included a grid search for learning rates ranging from 3e-6 to 3e-4, as well as an exploration of the hyperparameters most influential in each algorithm, such as $\beta$ in DPO.

In Table 7, we present the results of various algorithms, including the top-performing variants of DPO. The results indicate that RL can consistently outperform SFT in purely perceptual tasks, even when starting from a highly optimized pre-trained checkpoint, which is a non-trivial conclusion. As previously mentioned, purely perceptual tasks lack text-level reasoning, exploration, and reflection. While existing studies [66, 67, 68] have shown that reinforcement learning can enhance models that are either not pre-trained or insufficiently pre-trained, this improvement likely stems from mechanisms similar to those used by SFT. It remains uncertain whether reinforcement learning is beneficial for a model that has been pre-trained to its maximum potential.

In conjunction with the results of curriculum learning, we summarize the principles underlying the advantages of RL in post-training as follows, which are fundamentally different from those in RL work related to LLM reasoning:

- Due to the targeted nature of rollouts selection by these algorithms, they exhibit greater robustness in the data distribution during post-training. The advantages of RL in post-training may stem from the adaptive data distribution selection, akin to curriculum learning.

- Similarly, because DPO selects samples that include both correct and incorrect outputs, it may mitigate the impact of simplistic data and, more importantly, erroneous ground truth.

- Compared to SFT and curriculum learning, which use ground truth coordinates for training, RL relies solely on the model it-self's outputs. This results in a more gradual and stable training process.

## 5.5 Scaling

**Scaling settings.** We have set up the training settings for the scaling experiments by integrating all the conclusions drawn from previous sections. Specifically, we combined all available datasets at a certain ratio. Note that the Seeclick dataset was excluded after experimentation due to the excessive number of elements on average in each screenshot. For instance, in the case of 40M training, Table 8 shows the proportion of each dataset used. For other training volumes, such as 20M, the proportions are calculated accordingly based on the weights. During the training process, we further increased the batch size to 8K. As illustrated in Figure 12, the training loss curve is exceptionally smooth and stable. We also observed some fluctuations in the performance across different benchmarks during training, which led us to conduct tests at regular intervals and ultimately select the checkpoint that performed best across all benchmarks. Table 9 presents the model cards for all the models trained.

Table 8: Detailed training data proportion for Phi-Ground-4B-16C with 40M training volume. †: The number of samples here does not refer to the quantity of images or elements. In fact, because each element has three types of references—positional, appearance, and functional—we randomly combine them during training as model inputs. This combination could involve one, two, or all three types. Consequently, a single element can be paired with various references, resulting in multiple samples. This explains why the numbers here differ from those described earlier.

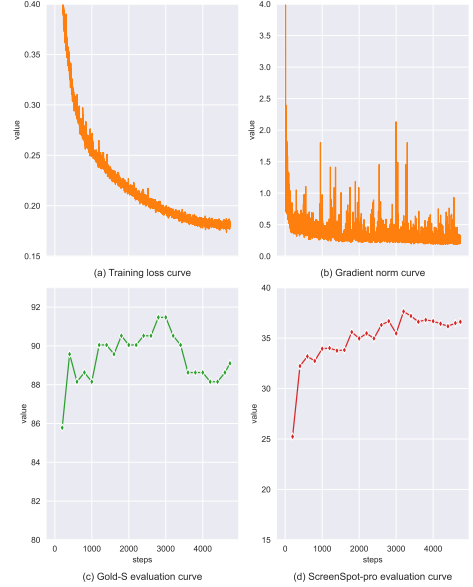| Dataset | Samples † | Epoch | Weight |
|---|---|---|---|
| BingSearch | 158 K | 7.6 | 3.0% |
| Linux | 149 K | 3.6 | 1.3% |
| MacOS | 69 K | 3.9 | 0.7% |
| Windows | 3.5 M | 1.1 | 10.0% |
| GUIAct | 155 K | 3.4 | 1.3% |
| E2ISynth | 180 K | 2.9 | 1.3% |
| Fineweb | 9.0 M | 0.9 | 21.0 % |
| CommonCrawl | 25.3 M | 1.0 | 60.0% |
| Human | 160 K | 3.5 | 1.4% |



Figure 12: Training and evaluation curves for Phi-Ground-4B-16C.

**Scaling effect of parameters-computation trade-off.** Scaling laws [38, 39, 40] typically examine the relationship between model parameters and training data given a fixed training budget. In our scenario, we further consider the computational cost during testing, which is not only related to the number of parameters but is also directly influenced by the number of image tokens. In perception tasks, the number of image tokens is highly correlated with the model's capability (see Table 2). In this section, we aim to investigate the relationship among model parameters, the number of image tokens, and training data volume. This is particularly relevant in real-world applications where developers are concerned with both training performance under limited resources and application latency (rather than just parameter count). Conducting such experiments can offer more economical strategies.

Table 9: Models card of Phi-Ground family.

| Name | Model details | Training details | Image tokens |
|---|---|---|---|
| Phi-Ground-4B-7C | Pretrained with Phi-3.5-Vision-Instruct as base model, which is a 4.1B VLM with CLIP as image encoder. During training we set the num_crops of the model to 6 and make sure the image's resolution is $2 \times 3$. The model will devide the image into $2 \times 3 + 1 = 7$ crops, in which there is a global crop. | Trained with 40M data; BS=8192; LR=8e-5. The training cost 200 A100 GPU days in total. | 1045 |
| Phi-Ground-4B-16C | Same with the above but change the image resolution to $3 \times 5$ and num_crops to 15. | Trained with 40M data; BS=8192; LR=8e-5. The training cost 440 A100 GPU days in total. | 2353 |
| Phi-Ground-4B-16C-DPO | Same with Phi-Ground-4B-16C | Three rounds of DPO finetuning with more desktop data from Phi-Ground-4B-16C. See Sec. 5.4 and Table 11 for more details. | 2353 |
| Phi-Ground-4B-29C | Same with Phi-Ground-4B-16C except for changing the image resolution to $4 \times 7$ and num_crops to 28. | Trained with 20M data; BS=8192; LR=8e-5. The training cost 415 A100 GPU days in total. | 4237 |
| Phi-Ground-7B-7C | A variation of Phi-4-MM [79] with 7B parameters. The model use SigLip as the image encoder. The input image's resolution is $2 \times 3$ and num_crops of the model is 6 . | Trained with 30M data; BS=8192; LR=1e-5. The training cost 350 A100 GPU days in total. | 1841 |
| Phi-Ground-7B-16C | Same with Phi-Ground-7B-7C except for changing the image resolution to $3 \times 5$ and num_crops to 15. | Trained with 15M data; BS=8192; LR=1e-5. The training cost 450 A100 GPU days in total. | 4161 |
| Phi-Ground-7B-16C-DPO | Same with Phi-Ground-7B-16C. | Three rounds of DPO finetuning with more desktop data from Phi-Ground-7B-16C. See Sec. 5.4 and Table 11 for more details. | 4161 |
| Phi-Ground-7B-29C | Same with Phi-Ground-7B-7C except for changing the image resolution to $4 \times 7$ and num_crops to 28. | Trained with 7.5M data; BS=8192; LR=1e-5. The training cost 390 A100 GPU days in total. | 7505 |

Specifically, we used the 4.1B Phi-3.5-Vision-Instruct model and the 7B Phi-4-MM model as the base models for training. For each model, we configured three different image settings corresponding to varying numbers of image

tokens. Specifically, the phi model family scales an image and crops it into a grid-shaped square, such as $336 \times 336$ for Phi-3.5-V. By configuring the model's num_crops parameter, we can adjust the maximum number of patches into which an image can be cropped. In our three image settings, we pad the images to shapes of $3 \times 2$, $5 \times 3$, and $7 \times 4$ using white padding, and set num_crops to 6, 15, and 28, respectively. Under these three configurations, the final images are divided into 7, 16, and 29 patches. This is because the model also includes a default global image patch.

During training, we fixed the total training budget for all models to 450 ($\pm 50$) NVIDIA A100-80G GPU days. Table 9 provides more detailed information on the image tokens and actual training durations for each model. In this manner, we trained six models under a fixed training budget. Our aim is to understand which model design can most effectively achieve optimal performance given the available computational resources.



Figure 13: Illustration of the evaluation results in relation to the training computation load. The Y-axis represents the benchmark scores in click accuracy, while the X-axis denotes the training computation per sample in TFLOPs. This training computation is estimated using the formula $FLOPs = 6ND$, where $N$ is the number of image tokens and $D$ is the number of model parameters.

Figure 13 presents the training results of these six models. Based on our evaluations, the inference time generally aligns with the relationship depicted on the x-axis of the graph. Many current studies typically report only the number of parameters when discussing model performance, without emphasizing computational aspects, such as the number of image tokens. In our experiments, where the model architecture is fixed, we observe that for more advanced and challenging benchmarks like ScreenSpot-pro and UI-Vision, the impact of image tokens is significant. Specifically, when the number of image tokens is low, it may become a bottleneck, resulting in the inability to perceive small objects and thus reducing the score. When the number of image tokens exceeds 2000, their impact gradually diminishes, meaning further increases in image tokens do not yield substantial marginal benefits akin to scaling laws.

In some test datasets that do not require high resolution, such as ScreenSpot-V2, neither the model size nor the number of image tokens significantly affects performance. Furthermore, when the number of image tokens exceeds the benchmark requirements (as previously mentioned bottleneck), the impact on perception is very limited, as illustrated by the results in the Figure 13 with o4-mini as the planner. However, the difference between using a planner and not using one is quite significant.

**Scaling post-training** We attempted to extend the duration of DPO post-training to observe its effects. Our findings indicate that conducting long-epoch Offline DPO training directly can lead to an initial increase in model performance followed by a decline, as illustrated by the gray line in Figure 14. We believe this phenomenon may be related to distribution shifting [80, 81], a concept frequently discussed in previous research. Consequently, we plan to update the rollouts more frequently. Given our stringent rollout selection criteria, we opted to conduct multiple rounds of DPO to approximate the effects of Online DPO [80] efficiently. Specifically, we performed a new rollout every 100 steps and initiated a new round of training. Ultimately, after three rounds of DPO, we achieved the best results, as shown in Figure 14.

Figure 14: Multi-turns DPO vs. Offline DPO for in-domain post-training.

In the post-training phase, as previously mentioned, we increased the proportion of in-domain data, which primarily consists of Web Search and human-labeled data. After training, the results for in-domain data (our Gold dataset) improved significantly. Surprisingly, several general benchmarks, such as ScreenSpot-Pro and UI-Vision, also showed noticeable improvements. This could be partly because these benchmarks have some overlap with our target software (e.g., PPT), and partly because the distribution of desktop applications is similar to that of these benchmarks. Other benchmarks not shown here maintained their original scores or experienced slight improvements after DPO, with detailed results available in the experimental tables in Appendix B.

## 6 Evaluation and Case Study

### 6.1 UI Grounding Benchmark Results
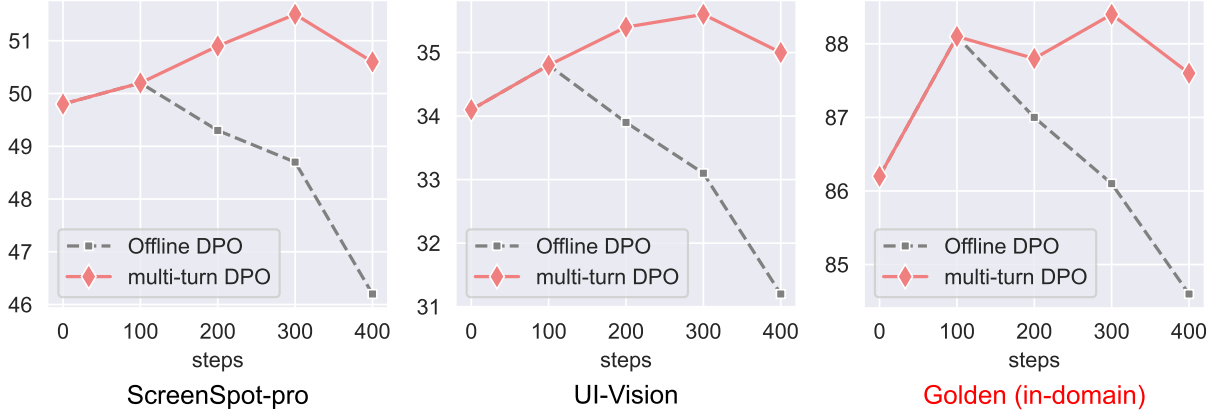
Table 10 presents the test results of several open-source models with fewer than $10B$ parameters on our selected five GUI grounding test sets. The results in the upper block were obtained using the benchmark's built-in reference expressions, typically an instruction or short REs. In contrast, the lower block shows the results when we used o4-mini to generate Long REs, which were then tested by the grounding models.

Table 10: The comparison of results across five GUI grounding test sets, which were tested by us, is presented. More detailed tables of results for additional open-source and closed-source models can be found in Appendix B.

| model | ScreenSpot-V2 | | | | ScreenSpot-pro | UI-Vision | | | | ShowDown | Gold | |
| | Desktop | Web | Mobile | AVG | AVG | basic | functional | spatial | AVG | AVG | Gold-S | ALL |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *End-to-end model setting (Use short REs)* | | | | | | | | | | | | |
| SeeClick-9.6B [26] | 64.6 | 49.7 | 43.9 | 55.1 | 1.1 | 9.4 | 4.7 | 2.1 | 5.4 | 24.6 | 51.7 | 20.4 |
| UGround-7B [30] | 73.2 | 78.3 | 72.7 | 76.1 | 16.5 | 11.5 | 12.2 | 2.8 | 8.8 | 46.5 | 74.9 | 54.9 |
| UGround-v1-7B [30] | 87.1 | 86.1 | 89.4 | 87.7 | 31.1 | 15.4 | 17.1 | 6.25 | 12.9 | 57.8 | 84.4 | 66.4 |
| OS-Atlas-4B [31] | 73.5 | 59.6 | 74.5 | 71.9 | 3.7 | - | - | - | - | 15.8 | 47.9 | 22.0 |
| OS-Atlas-7B [31] | 85.5 | 77.2 | 84.0 | 84.1 | 18.9 | 12.2 | 11.2 | 3.67 | 9.0 | 41.1 | 66.4 | 48.8 |
| UI-TARS-2B [28] | 87.2 | 79.7 | 82.8 | 84.7 | 27.7 | - | - | - | - | 59.8 | 79.2 | 60.1 |
| UI-TARS-7B [28] | **93.0** | **90.2** | 89.4 | **91.6** | 35.7 | 20.1 | 24.3 | 8.4 | 17.6 | 66.1 | **87.2** | 76.8 |
| UI-TARS-1.5-7B [29] | 86.9 | 87.6 | **90.0** | 89.0 | 42.6 | 28.8 | 27.5 | **10.7** | 22.3 | **67.2** | 86.7 | 77.2 |
| **Phi-Ground-4B-16C-DPO** | 84.4 | 86.4 | 78.1 | 83.1 | 38.0 | 33.4 | 34.3 | 5.8 | 24.5 | 58.2 | **87.2** | 78.2 |
| **Phi-Ground-7B-16C-DPO** | 83.3 | 84.8 | 79.3 | 83.8 | **43.2** | **36.8** | **37.1** | 7.6 | **27.2** | 62.5 | 84.4 | **79.6** |
| *Agent setting (Use long REs) with O4-mini as planner* | | | | | | | | | | | | |
| SeeClick-9.6B | 43.5 | 64.4 | 43.3 | 53.5 | 1.2 | 5.2 | 5.0 | 3.2 | 4.5 | 19.6 | 39.4 | 15.6 |
| UGround-7B | 84.6 | 90.0 | 88.5 | 89.0 | 23.7 | 25.1 | 22.8 | 11.3 | 19.7 | 62.4 | 83.2 | 59.5 |
| UGround-v1-7B | 89.8 | 92.2 | 91.9 | 92.1 | 32.5 | 30.5 | 29.3 | 13.9 | 24.6 | 66.7 | 88.0 | 73.8 |
| OS-Atlas-4B | 41.9 | 64.6 | 57.7 | 57.4 | 2.0 | - | - | - | - | 18.0 | 37.0 | 18.2 |
| OS-Atlas-7B | 80.6 | 86.2 | 88.9 | 84.7 | 21.1 | 17.2 | 16.6 | 7.5 | 13.8 | 45.5 | 68.3 | 52.2 |
| UI-TARS-2B | 85.5 | 92.2 | 90.0 | 90.4 | 35.9 | - | - | - | - | 66.4 | 85.6 | 77.1 |
| UI-TARS-7B | 91.5 | 93.1 | 92.9 | 93.0 | 40.6 | 33.2 | 33.4 | 16.7 | 27.8 | 69.8 | 90.9 | 81.5 |
| UI-TARS-1.5-7B | 89.9 | **92.7** | 92.0 | 92.2 | 48.8 | 35.1 | 35.1 | 17.9 | 29.4 | 71.6 | 90.4 | 80.3 |
| **Phi-Ground-4B-16C-DPO** | **92.8** | **92.7** | 92.4 | 92.3 | 51.5 | 43.8 | 42.1 | **21.0** | 35.6 | 73.5 | **95.2** | **88.4** |
| **Phi-Ground-7B-16C-DPO** | 92.6 | 92.6 | **93.0** | **93.4** | **55.0** | **44.2** | **43.8** | 20.5 | **36.2** | **73.9** | 93.8 | 88.2 |

Our grounding model is trained specifically for the agent setting, meaning that the training dataset primarily consists of various combinations of REs. As a result, our model achieves significant advantages and SOTA results across all benchmarks in the agent setting. Specifically, ScreenSpot-pro achieved an accuracy of 55.0. UI-Vision also attained a result of 36.2, which is the highest for this benchmark. Furthermore, our results on the Showdown benchmark surpass those of commercial models like OpenAI Operator and Claude Computer Use (see Table 15 in the Appendix).

In the end-to-end model setting, our model consistently outperforms others on the ScreenSpot-Pro, UI-Vision, and our Gold dataset, although its performance on ScreenSpot-V2 is relatively average. ScreenSpot V1 and V2 (with V2 sharing most of its data with V1) have long been crucial benchmarks for GUI grounding testing. We observed that some models perform well on ScreenSpot-V2 but do not demonstrate the same significant advantages on newly emerging benchmarks. This could be a result of developers optimizing their models based on a single benchmark over time. In our approach, we did not include any mobile data in the training set (as this is not our focus scenario), nor did we balance the training set with icon and text-based buttons (since the scenarios faced by product users are not balanced either). However, these techniques might significantly impact the accuracy on ScreenSpot-V1 and V2. Throughout our development process, the selection and ablation of each technique were carefully considered across multiple benchmarks. As a result, our model exhibits more balanced performance and better generalization.

## 6.2 Error Analysis



Figure 15: Types and Proportions of Errors on the ScreenSpot-pro Benchmark. In each image, the red rectangles represent the regions corresponding to the ground truth. Red circles indicate erroneous outputs from the previous stage, while green circles denote correct outputs from the current stage. The centers of the green circles fall within the ground truth boundaries. To avoid obstructing the image content, we have enlarged the green circles in some of the images.

To analyze and illustrate the errors made by current grounding models as a case study, we selected a challenging benchmark, ScreenSpot-Pro, as an example. We employed the Phi-Ground-4B-DPO as the grounding model and designed a cascading approach to sequentially process the benchmark data.

Specifically, as shown in Figure 15, **Stage A** involves using the benchmark-provided instructions as reference expressions (also known as short REs) for the total of 1,534 test samples from ScreenSpot-Pro. As previously mentioned, our model successfully resolved 38% of the test cases, leaving 951 incorrect samples. In **Stage B**, we used O4-mini as a planner for the remaining 951 samples to generate long REs as input for the model. This approach further resolved 292 samples, leaving 659 samples unresolved. Across both stages, we successfully addressed 57% of the benchmark samples. In **Stage C**, for the remaining 659 samples, we intended to use human experts as planners to generate REs to observe the potential error rates in the planner section. However, this approach was deemed too costly. Therefore, we first had O4-mini generate REs using the Long-Gold method (see Sec. 4.2), which involves disclosing the GT bbox during RE generation. Human experts then reviewed the samples, GT, and REs generated by O4-mini to correct any errors and

produce the final REs for model input. This stage resolved an additional 236 samples. In **Stage D**, we recognized that for REs of similar quality, different styles and emphases might lead to varying results. Thus, we repeated the RE and grounding process of Stage C seven times for the remaining 423 samples, resolving an additional 88 samples. Ultimately, 335 samples remained with errors largely unrelated to RE quality. We will analyze the errors and their implications in each stage in detail below.

We first observe that end-to-end grounding models lack spatial reasoning capabilities, as illustrated in the example from Stage B of Figure 15. When certain keywords appear in the instruction, such as "keyword" in Example 1 or "screen" in Example 2, the grounding model tends to directly highlight the locations of these words in the image. However, in Example 1, the interactive region is actually the white rectangular input box, and clicking on the text of the label might result in a failed interaction. Such spatial reasoning requires a degree of common sense, rather than being purely a grounding task. The introduction of a planner addresses this type of task effectively by directly describing "white rectangle" in the RE. We refer to such errors as "**planning omissions**," which account for 19% of the total sample and 30.7% of the total errors.

However, the planning of O4-mini may also encounter errors, particularly in scenarios where the target area contains multiple similar regions or when specialized application knowledge is required. In such cases, the planner's hallucinations can lead to mistakes. For instance, as illustrated in Figure 15 Stage C, the markdown display button is typically located in the upper right corner. However, in the example shown, two work pages are open, causing the button to be centered. This resulted in an incorrect RE by O4-mini, leading to erroneous grounding. After manually correcting the erroneous RE, our model was able to produce the correct result. We refer to this type of error as a "**planning error**," which accounts for 15.4% of the total samples and 24.8% of the total errors.

For the remaining samples, we observed that the grounding model might exhibit a preference for a certain style of RE, even when the quality is consistent, as evidenced by the pass@8 metric indicating new correct samples. However, this influence is minor, affecting only 5.7% of the samples and accounting for 9.3% of the errors. Additionally, we found that 117 samples (7.6% of the samples and 12.3% of the errors) were impacted because the target area or its vicinity contained languages not covered by our model, such as Chinese. Due to the strict selection of training data, primarily from the CommonCrawl dataset, which exclusively includes only English data, the model failed to correctly recognize many straightforward and easy situations due to language issues. Ultimately, 218 samples remained as particularly challenging data, making up another 14.2% of the samples and 22.9% of the errors.

For the remaining errors, we provide a more detailed case study in Appendix F. We categorize these errors into several identifiable types:

- Accuracy issues arise due to excessively extreme screen sizes and shapes. For instance, screens with an ultra-wide aspect ratio may result in output coordinates that deviate from the intended target area.
- Language descriptions fail to adequately constrain areas of spatial planning, such as when generating tables and instructing to click on a blank cell in the 13th row and 8th column. Such specific spatial tasks present significant challenges for both the planner and the grounding model.
- Regions that are difficult to describe using natural language.

# 7 Social Impacts and Open Questions

With the development of CUA, we have both expectations and concerns regarding this direction. Primarily, there is the issue of user privacy. During our training process, we have verified the legality of the licenses for the open-source datasets used, and ensured that licenses are valid in Bing search filtering and web filtering. However, when CUA is successfully deployed in user environments in the future, the need for grounding and planning may require screenshots of users' screens to be uploaded to the cloud, potentially leading to privacy breaches. Throughout the entire research and product deployment process, we may need to establish relevant protocols, legal frameworks, or algorithms to ensure the protection of user privacy.

Secondly, there is the issue of accountability for erroneous actions performed by CUAs. There are instances where CUAs might execute irreversible and harmful operations, such as closing software without saving files or even deleting important documents. At the system level, we need to explore human-computer collaboration methods that allow CUAs to efficiently replace human labor while ensuring human oversight. From the perspective of GUI grounding, we have observed that errors due to incorrect grounding can have more severe consequences. This is because a trained grounding model, when making mistakes, still outputs interactively meaningful regions rather than blank areas, thus increasing the likelihood of irreversible impacts. For instance, the multiplication symbol on a calculator might be mistakenly interpreted as a command to close software due to similar symbols, leading to unintended software closure. Some recent studies [82, 83] have attempted to use MLLM to verify actions post-grounding, but these have shown limited

effectiveness and increased time costs. Developing a benchmark to evaluate the potential harmfulness of GUI grounding models could also be highly beneficial.

## 8 Conclusion

In conclusion, we have developed the Phi-Ground model family, which significantly enhances GUI grounding capabilities by improving the perception of interactive elements in digital interfaces. Our comprehensive empirical study identified critical factors such as data distribution, input/output formats, and computational efficiency that influence model performance. Using a two-stage approach, we combined advanced MLLMs for generating detailed REs with a specialized grounding model for precise coordinate output, achieving state-of-the-art results across various benchmarks, including challenging ones like ScreenSpot-pro and UI-Vision. While our models demonstrate promising results, we acknowledge the societal implications of deploying CUAs, especially regarding user privacy and error accountability. Our research not only advances GUI grounding but also offers insights applicable to other multimodal perception tasks, contributing to the development of more reliable and efficient CUAs.

## References

[1] Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 2024.

[2] Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. The rise and potential of large language model based agents: A survey. *Science China Information Sciences*, 2025.

[3] Yuheng Cheng, Ceyao Zhang, Zhengwen Zhang, Xiangrui Meng, Sirui Hong, Wenhao Li, Zihao Wang, Zekai Wang, Feng Yin, Junhua Zhao, et al. Exploring large language model based intelligent agents: Definitions, methods, and prospects. *arXiv preprint arXiv:2401.03428*, 2024.

[4] Chaoyun Zhang, Shilin He, Jiaxu Qian, Bowen Li, Liqun Li, Si Qin, Yu Kang, Minghua Ma, Guyue Liu, Qingwei Lin, et al. Large language model-brained gui agents: A survey. *arXiv preprint arXiv:2411.18279*, 2024.

[5] Pascal J Sager, Benjamin Meyer, Peng Yan, Rebekka von Wartburg-Kottler, Layan Etaiwi, Aref Enayati, Gabriel Nobel, Ahmed Abdulkadir, Benjamin F Grewe, and Thilo Stadelmann. A comprehensive survey of agents for computer use: Foundations, challenges, and future directions. *arXiv preprint arXiv:2501.16150*, 2025.

[6] Mohsen Soori, Behrooz Arezoo, and Roza Dastres. Artificial intelligence, machine learning and deep learning in advanced robotics, a review. *Cognitive Robotics*, 2023.

[7] Demetris Vrontis, Michael Christofi, Vijay Pereira, Shlomo Tarba, Anna Makrides, and Eleni Trichina. Artificial intelligence, robotics, advanced technologies and human resource management: a systematic review. *Artificial intelligence and international HRM*, 2023.

[8] OpenAI. Introducing openai o3 and o4-mini, 2025.

[9] Anthropic. Claude sonnet 4, 2025.

[10] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

[11] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.

[12] OpenAI. Operator system card, 2025.

[13] Anthropic. Introducing computer use, a new claude 3.5 sonnet, and claude 3.5 haiku, 2024.

[14] Michael S Greenberg, Jennifer C Byington, and David G Harper. Mobile agents and security. *IEEE Communications magazine*, 1998.

[15] Jose M Such, Agustín Espinosa, and Ana García-Fornes. A survey of privacy in multi-agent systems. *The Knowledge Engineering Review*, 2014.

[16] Sohye Lim and Hongjin Shim. No secrets between the two of us: Privacy concerns over using ai agents. *Cyberpsychology: Journal of Psychosocial Research on Cyberspace*, 2022.

[17] K Cartrysse and JCA Van Der Lubbe. Privacy in mobile agents. In *IEEE First Symposium onMulti-Agent Security and Survivability, 2004*. IEEE, 2004.

[18] Evan Zheran Liu, Kelvin Guu, Panupong Pasupat, Tianlin Shi, and Percy Liang. Reinforcement learning on web interfaces using workflow-guided exploration. *ICLR*, 2018.

[19] Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su. Mind2web: Towards a generalist agent for the web. *NIPS*, 2023.

[20] Yiheng Xu, Zekun Wang, Junli Wang, Dunjie Lu, Tianbao Xie, Amrita Saha, Doyen Sahoo, Tao Yu, and Caiming Xiong. Aguvis: Unified pure vision agents for autonomous gui interaction. *ICML*, 2025.

[21] Wentong Chen, Junbo Cui, Jinyi Hu, Yujia Qin, Junjie Fang, Yue Zhao, Chongyi Wang, Jun Liu, Guirong Chen, Yupeng Huo, et al. Guicourse: From general vision language models to versatile gui agents. *arXiv preprint arXiv:2406.11317*, 2024.

[22] Segev Shlomov, Aviad Sela, Ido Levy, Liane Galanti, Roy Abitbol, et al. From grounding to planning: Benchmarking bottlenecks in web agents. *arXiv preprint arXiv:2409.01927*, 2024.

[23] Suyu Ye, Haojun Shi, Darren Shih, Hyokun Yun, Tanya Roosta, and Tianmin Shu. Realwebassist: A benchmark for long-horizon web assistance with real-world users. *arXiv preprint arXiv:2504.10445*, 2025.

[24] Shravan Nayak, Xiangru Jian, Kevin Qinghong Lin, Juan A Rodriguez, Montek Kalsi, Rabiul Awal, Nicolas Chapados, M Tamer Özsu, Aishwarya Agrawal, David Vazquez, et al. Ui-vision: A desktop-centric gui benchmark for visual perception and interaction. *arXiv preprint arXiv:2503.15661*, 2025.

[25] Kaixin Li, Ziyang Meng, Hongzhan Lin, Ziyang Luo, Yuchen Tian, Jing Ma, Zhiyong Huang, and Tat-Seng Chua. Screenspot-pro: Gui grounding for professional high-resolution computer use. *arXiv preprint arXiv:2504.07981*, 2025.

[26] Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Li YanTao, Jianbing Zhang, and Zhiyong Wu. Seeclick: Harnessing gui grounding for advanced visual gui agents. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2024.

[27] Yuhang Liu, Pengxiang Li, Congkai Xie, Xavier Hu, Xiaotian Han, Shengyu Zhang, Hongxia Yang, and Fei Wu. Infigui-r1: Advancing multimodal gui agents from reactive actors to deliberative reasoners. *arXiv preprint arXiv:2504.14239*, 2025.

[28] Yujia Qin, Yining Ye, Junjie Fang, Haoming Wang, Shihao Liang, Shizuo Tian, Junda Zhang, Jiahao Li, Yunxin Li, Shijue Huang, et al. Ui-tars: Pioneering automated gui interaction with native agents. *arXiv preprint arXiv:2501.12326*, 2025.

[29] ByteDance Seed. Ui-tars-1.5. `https://seed-tars.com/1.5`, 2025.

[30] Boyu Gou, Ruohan Wang, Boyuan Zheng, Yanan Xie, Cheng Chang, Yiheng Shu, Huan Sun, and Yu Su. Navigating the digital world as humans do: Universal visual grounding for gui agents. *ICLR*, 2025.

[31] Zhiyong Wu, Zhenyu Wu, Fangzhi Xu, Yian Wang, Qiushi Sun, Chengyou Jia, Kanzhi Cheng, Zichen Ding, Liheng Chen, Paul Pu Liang, et al. Os-atlas: A foundation action model for generalist gui agents. *ICLR*, 2025.

[32] Zhengyuan Yang, Zhe Gan, Jianfeng Wang, Xiaowei Hu, Faisal Ahmed, Zicheng Liu, Yumao Lu, and Lijuan Wang. Unitab: Unifying text and box outputs for grounded vision-language modeling. In *ECCV*. Springer, 2022.

[33] Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. *NIPS*, 2023.

[34] Haoxuan You, Haotian Zhang, Zhe Gan, Xianzhi Du, Bowen Zhang, Zirui Wang, Liangliang Cao, Shih-Fu Chang, and Yinfei Yang. Ferret: Refer and ground anything anywhere at any granularity. *ICLR*, 2024.

[35] Jiasen Lu, Christopher Clark, Rowan Zellers, Roozbeh Mottaghi, and Aniruddha Kembhavi. Unified-io: A unified model for vision, language, and multi-modal tasks. *ICLR*, 2023.

[36] Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu, and Lijuan Wang. Git: A generative image-to-text transformer for vision and language. *arXiv preprint arXiv:2205.14100*, 2022.

[37] Parvinder Kaur, Baljit Singh Khehra, and Er Bhupinder Singh Mavi. Data augmentation for object detection: A review. In *2021 IEEE International Midwest Symposium on Circuits and Systems (MWSCAS)*. IEEE, 2021.

[38] Joel Hestness, Sharan Narang, Newsha Ardalani, Gregory Diamos, Heewoo Jun, Hassan Kianinejad, Md Mostofa Ali Patwary, Yang Yang, and Yanqi Zhou. Deep learning scaling is predictable, empirically. *arXiv preprint arXiv:1712.00409*, 2017.

[39] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.

[40] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. In *NIPS*, 2022.

[41] Xinyi Liu, Xiaoyi Zhang, Ziyun Zhang, and Yan Lu. Ui-e2i-synth: Advancing gui grounding with large-scale instruction synthesis. *arXiv preprint arXiv:2504.11257*, 2025.

[42] Hugh Zhang, Jeff Da, Dean Lee, Vaughn Robinson, Catherine Wu, William Song, Tiffany Zhao, Pranav Raja, Charlotte Zhuang, Dylan Slack, et al. A careful examination of large language model performance on grade school arithmetic. *NIPS*, 2024.

[43] Tianwen Wei, Liang Zhao, Lichang Zhang, Bo Zhu, Lijie Wang, Haihua Yang, Biye Li, Cheng Cheng, Weiwei Lü, Rui Hu, et al. Skywork: A more open bilingual foundation model. *arXiv preprint arXiv:2310.19341*, 2023.

[44] General Agents Team. The showdown computer control evaluation suite, 2025.

[45] Zhipeng Huang, Zhizheng Zhang, Yiting Lu, Zheng-Jun Zha, Zhibo Chen, and Baining Guo. Visualcritic: Making lmms perceive visual quality like humans. *arXiv preprint arXiv:2403.12806*, 2024.

[46] Rizhao Cai, Zirui Song, Dayan Guan, Zhenhao Chen, Yaohang Li, Xing Luo, Chenyu Yi, and Alex Kot. Benchlmm: Benchmarking cross-style visual capability of large multimodal models. In *ECCV*. Springer, 2024.

[47] Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A Smith, Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but not perceive. In *ECCV*. Springer, 2024.

[48] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multimodal llm's referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023.

[49] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. Ofa: Unifying architectures, tasks, and modalities through a simple sequence-to-sequence learning framework. In *ICML*. PMLR, 2022.

[50] Common crawl - open repository of web crawl data, 2025.

[51] Yadong Lu, Jianwei Yang, Yelong Shen, and Ahmed Awadallah. Omniparser for pure vision based gui agent. *arXiv preprint arXiv:2408.00203*, 2024.

[52] Zheng Hui, Yinheng Li, Tianyi Chen, Colby Banbury, Kazuhito Koishida, et al. Winclick: Gui grounding with multimodal large language models. *arXiv preprint arXiv:2503.04730*, 2025.

[53] OpenAI. Hello gpt-4o, 2024.

[54] Samuel Lavoie, Polina Kirichenko, Mark Ibrahim, Mahmoud Assran, Andrew Gordon Wilson, Aaron Courville, and Nicolas Ballas. Modeling caption diversity in contrastive vision-language pretraining. *arXiv preprint arXiv:2405.00740*, 2024.

[55] Ziqiang Xu, Qi Dai, Tian Xie, Yifan Yang, Kai Qiu, DongDong Chen, Zuxuan Wu, and Chong Luo. Viarl: Adaptive temporal grounding via visual iterated amplification reinforcement learning. *arXiv preprint arXiv:2505.15447*, 2025.

[56] Sara Ghazanfari, Alexandre Araujo, Prashanth Krishnamurthy, Siddharth Garg, and Farshad Khorrami. Emma: Efficient visual alignment in multi-modal llms. *arXiv preprint arXiv:2410.02080*, 2024.

[57] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

[58] Walter R Gilks and Pascal Wild. Adaptive rejection sampling for gibbs sampling. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 41(2), 1992.

[59] Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.

[60] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.

[61] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *NIPS*, 2023.

[62] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3), 1992.

[63] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

[64] Jian Hu. Reinforce++: A simple and efficient approach for aligning large language models. *arXiv preprint arXiv:2501.03262*, 2025.

[65] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Yang Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.

[66] Zhengxi Lu, Yuxiang Chai, Yaxuan Guo, Xi Yin, Liang Liu, Hao Wang, Han Xiao, Shuai Ren, Guanjing Xiong, and Hongsheng Li. Ui-r1: Enhancing efficient action prediction of gui agents by reinforcement learning. *arXiv preprint arXiv:2503.21620*, 2025.

[67] Yuqi Zhou, Sunhao Dai, Shuai Wang, Kaiwen Zhou, Qinglin Jia, and Jun Xu. Gui-g1: Understanding r1-zero-like training for visual grounding in gui agents. *arXiv preprint arXiv:2505.15810*, 2025.

[68] Run Luo, Lu Wang, Wanwei He, and Xiaobo Xia. Gui-r1: A generalist r1-style vision-language action model for gui agents. *arXiv preprint arXiv:2504.10458*, 2025.

[69] Leandro von Werra, Younes Belkada, Lewis Tunstall, Edward Beeching, Tristan Thrush, Nathan Lambert, Shengyi Huang, Kashif Rasul, and Quentin Gallouédec. Trl: Transformer reinforcement learning. `https://github.com/huggingface/trl`, 2020.

[70] Yao Zhao, Rishabh Joshi, Tianqi Liu, Misha Khalman, Mohammad Saleh, and Peter J Liu. Slic-hf: Sequence likelihood calibration with human feedback. *arXiv preprint arXiv:2305.10425*, 2023.

[71] Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal Valko, and Daniele Calandriello. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 2024.

[72] Haozhe Ji, Cheng Lu, Yilin Niu, Pei Ke, Hongning Wang, Jun Zhu, Jie Tang, and Minlie Huang. Towards efficient exact optimization of language model alignment. In *ICML*, 2024.

[73] Huayu Chen, Guande He, Lifan Yuan, Ganqu Cui, Hang Su, and Jun Zhu. Noise contrastive alignment of language models with explicit rewards. *Advances in Neural Information Processing Systems*, 37:117784–117812, 2024.

[74] Sayak Ray Chowdhury, Anush Kini, and Nagarajan Natarajan. Provably robust dpo: Aligning language models with noisy feedback. In *International Conference on Machine Learning*, pages 42258–42274. PMLR, 2024.

[75] Yue Wu, Zhiqing Sun, Huizhuo Yuan, Kaixuan Ji, Yiming Yang, and Quanquan Gu. Self-play preference optimization for language model alignment. *arXiv preprint arXiv:2405.00675*, 2024.

[76] Igor Melnyk, Youssef Mroueh, Brian Belgodere, Mattia Rigotti, Apoorva Nitsure, Mikhail Yurochkin, Kristjan Greenewald, Jiri Navratil, and Jarret Ross. Distributional preference alignment of llms via optimal transport. *NIPS*, 2024.

[77] Chris Lu, Samuel Holt, Claudio Fanconi, Alex Chan, Jakob Foerster, Mihaela van der Schaar, and Robert Lange. Discovering preference optimization algorithms with and for large language models. *NIPS*, 2024.

[78] Karel D'Oosterlinck, Winnie Xu, Chris Develder, Thomas Demeester, Amanpreet Singh, Christopher Potts, Douwe Kiela, and Shikib Mehri. Anchored preference optimization and contrastive revisions: Addressing underspecification in alignment. *ACL*, 2025.

[79] Abdelrahman Abouelenin, Atabak Ashfaq, Adam Atkinson, Hany Awadalla, Nguyen Bach, Jianmin Bao, Alon Benhaim, Martin Cai, Vishrav Chaudhary, Congcong Chen, et al. Phi-4-mini technical report: Compact yet powerful multimodal language models via mixture-of-loras. *arXiv preprint arXiv:2503.01743*, 2025.

[80] Biqing Qi, Pengfei Li, Fangyuan Li, Junqi Gao, Kaiyan Zhang, and Bowen Zhou. Online dpo: Online direct preference optimization with fast-slow chasing. *arXiv preprint arXiv:2406.05534*, 2024.

[81] Yi Ren and Danica J Sutherland. Learning dynamics of llm finetuning. *ICLR*, 2024.

[82] Tiange Luo, Lajanugen Logeswaran, Justin Johnson, and Honglak Lee. Visual test-time scaling for gui agent grounding. *arXiv preprint arXiv:2505.00684*, 2025.

[83] Jungjae Lee, Dongjae Lee, Chihun Choi, Youngmin Im, Jaeyoung Wi, Kihong Heo, Sangeun Oh, Sunjae Lee, and Insik Shin. Safeguarding mobile gui agent via logic-based action verification. *arXiv preprint arXiv:2503.18492*, 2025.

# A Experiment Settings

Due to resource constraints and the evolution of the development process, different ablation experiments may have utilized varying hyper-parameters and data configurations. We have documented the detailed setup for each ablation experiment in the table below.

Table 11: Detailed training data configuration and hyper-parameters.

| EXP | Data configure | Hyper parameters |
|---|---|---|
| Tab. 2, Tab. 3 | Training samples: 2M<br>Dataset: [(SeeClick, 0.075 ep), (Fineweb, 0.062 ep),<br>(Windows, 0.088 ep), (MaxOS, 0.77 ep),<br>(Linux, 0.45 ep), (GUIAct, 0.44 ep)] | lr=8e-5; batch size=2048;<br>weight decay=0.01;<br>random-resize=0.5<br>max-grad-norm=0.1;<br>warmup-steps=50 |
| Tab. 4, Tab. 17 | Training samples: 5M<br>Dataset: [(SeeClick, 0.19 ep), (Fineweb, 0.16 ep),<br>(Windows, 0.22 ep), (MaxOS, 0.77 ep),<br>(Linux, 0.65 ep), (GUIAct, 0.63 ep)]<br>(BingSearch, 0.95 ep) | lr=8e-5; batch size=2048;<br>weight decay=0.01;<br>max-grad-norm=0.1;<br>warmup-steps=100 |
| Tab. 5 | Training samples: 5M<br>Base: [(SeeClick, 0.19 ep), (Fineweb, 0.16 ep),<br>(Windows, 0.23 ep), (MaxOS, 0.77 ep),<br>(Linux, 0.65 ep), (GUIAct, 0.63 ep)] | lr=8e-5; batch size=2048;<br>weight decay=0.01;<br>max-grad-norm=0.1;<br>warmup-steps=100 |
| Tab. 6 | Training samples: 5M<br>A-pretrain: [(SeeClick, 0.13 ep), (Fineweb, 0.11 ep),<br>(Windows, 0.15 ep), (MaxOS, 0.51 ep),<br>(Linux, 0.43 ep), (GUIAct, 0.52 ep),<br>(BingSearch, 0.95 ep), (CC, 0.06 ep)] | lr=8e-5; batch size=2048;<br>weight decay=0.01;<br>max-grad-norm=0.1;<br>warmup-steps=100 |
| Phi-Ground pre-train | In Table 8 and 9 | weight decay=0.01;<br>max-grad-norm=0.1;<br>warmup-ratio=6%;<br>Others in Table 9 |
| Phi-Ground DPO | Training samples: 400K for each round<br>DPO: [(Human, 0.5ep),<br>(Windows, 0.03 ep), (MaxOS, 0.14 ep),<br>(Linux, 0.2 ep), (GUIAct, 0.2 ep),<br>(BingSearch, 0.33 ep), (CC, 0.004 ep)] | lr=1e-5; batch size=256;<br>weight decay=0.01;<br>max-grad-norm=0.1;<br>warmup-steps=100 |

# B   Detailed Evaluation Results

Table 12: Detailed ScreenSpot-V2 results.

| Model | Desktop | | Web | | Mobile | | AVG. |
|---|---|---|---|---|---|---|---|
| | Text | Icon/Widget | Text | Icon/Widget | Text | Icon/Widget | |
| *End-to-end model setting (Use short REs)* | | | | | | | |
| SeeClick | 78.4 | 50.7 | 70.1 | 29.3 | 55.2 | 32.5 | 55.1 |
| UGround-7B[*] | 85.1 | 61.2 | 84.6 | 71.9 | 84.3 | 61.1 | 76.1 |
| UGround-v1-7B[*] | 83.6 | **90.5** | 85.8 | 86.3 | 95.5 | 83.2 | 87.7 |
| OS-Atlas-4B | 87.2 | 59.7 | 72.7 | 46.4 | 85.9 | 63.1 | 71.9 |
| OS-Atlas-7B | 95.2 | 75.8 | 90.7 | 63.6 | 90.6 | 77.3 | 84.1 |
| UI-TARS-2B | 95.2 | 79.1 | 90.7 | 68.6 | 87.2 | 78.3 | 84.7 |
| UI-TARS-7B | <u>96.9</u> | <u>89.1</u> | **95.4** | **85.0** | 93.6 | **85.2** | **91.6** |
| UI-TARS-1.5-7B[*] | 92.2 | 81.5 | 91.0 | <u>84.2</u> | <u>95.5</u> | <u>84.5</u> | <u>89.0</u> |
| **Phi-Ground-4B-7C** | 87.6 | 75.7 | 93.6 | 71.4 | 94.1 | 54.3 | 80.8 |
| **Phi-Ground-4B-16C** | 90.2 | 77.1 | 92.3 | 76.8 | 93.7 | 63.5 | 83.4 |
| **Phi-Ground-4B-16C-DPO** | 91.7 | 77.1 | <u>94.4</u> | 78.3 | 94.1 | 62.0 | 84.1 |
| **Phi-Ground-4B-29C** | 90.7 | 80.0 | 92.3 | 76.8 | 92.0 | 63.5 | 83.4 |
| **Phi-Ground-7B-7C** | 88.1 | 70.0 | 93.2 | 71.4 | 94.4 | 60.6 | 81.3 |
| **Phi-Ground-7B-16C** | 83.9 | 72.1 | 94.0 | 73.4 | 94.4 | 63.5 | 81.8 |
| **Phi-Ground-7B-16C-DPO** | 90.2 | 76.4 | 93.6 | 75.9 | **96.5** | 62.0 | 83.8 |
| **Phi-Ground-7B-29C** | **97.2** | 77.9 | 91.9 | 71.9 | 95.1 | 63.5 | 82.5 |
| *Agent setting (Use long REs) with GPT-4O as planner* | | | | | | | |
| SeeClick[*] | 55.7 | 28.8 | 43.6 | 24.1 | 67.9 | 55.3 | 48.1 |
| UGround-7B[*] | 90.2 | 68.3 | 91.9 | 72.9 | 95.8 | 76.4 | 84.3 |
| UGround-v1-7B[*] | 83.1 | **92.9** | 85.8 | 84.3 | **97.6** | **85.6** | <u>88.5</u> |
| OS-Atlas-4B[*] | 49.5 | 30.1 | 67.0 | 48.7 | 78.6 | 58.7 | 58.7 |
| OS-Atlas-7B[*] | 85.0 | 66.9 | 84.3 | 68.6 | 94.4 | 75.5 | 80.9 |
| UI-TARS-2B[*] | 91.2 | 77.7 | 90.6 | <u>77.3</u> | 95.8 | 81.7 | 86.9 |
| UI-TARS-7B[*] | **93.3** | <u>83.5</u> | 91.5 | **78.3** | 96.2 | <u>85.0</u> | **88.8** |
| UI-TARS-1.5-7B[*] | <u>92.8</u> | 80.6 | 91.5 | 75.4 | 96.5 | 84.1 | 87.8 |
| **Phi-Ground-4B-7C** | 92.7 | 76.2 | **94.6** | 75.4 | 95.8 | 76.0 | 85.5 |
| **Phi-Ground-4B-16C** | 90.2 | 80.0 | 92.7 | 72.9 | 96.5 | 76.9 | 86.0 |
| **Phi-Ground-4B-16C-DPO** | 91.2 | 80.0 | 92.7 | 73.4 | 97.2 | 76.9 | 86.4 |
| **Phi-Ground-4B-29C** | 90.2 | 79.3 | 92.7 | 74.5 | 96.2 | 76.4 | 86.1 |
| **Phi-Ground-7B-7C** | 90.2 | 80.7 | 92.3 | 74.4 | <u>97.6</u> | 78.8 | 86.8 |
| **Phi-Ground-7B-16C** | 89.6 | 81.4 | <u>93.6</u> | 73.4 | 96.5 | 79.3 | 86.7 |
| **Phi-Ground-7B-16C-DPO** | 92.2 | 80.0 | 93.2 | 73.9 | 92.2 | 80.0 | 87.2 |
| **Phi-Ground-7B-29C** | 91.2 | 80.7 | 93.2 | 75.4 | 97.2 | 77.9 | 87.0 |
| *Agent setting (Use long REs) with O4-mini as planner* | | | | | | | |
| SeeClick[*] | 60.6 | 26.3 | 76.8 | 51.9 | 57.1 | 29.5 | 53.5 |
| UGround-7B[*] | 92.6 | 76.6 | <u>97.5</u> | 82.5 | 92.9 | 84.0 | 89.0 |
| UGround-v1-7B[*] | 92.6 | 86.9 | **97.9** | 86.4 | 93.8 | 91.0 | 92.1 |
| OS-Atlas-4B[*] | 50.5 | 33.3 | 73.9 | 55.3 | 63.0 | 52.4 | 57.4 |
| OS-Atlas-7B[*] | 86.0 | 75.2 | 94.7 | 77.7 | 86.0 | 81.8 | 84.7 |
| UI-TARS-2B[*] | 93.6 | 77.4 | **97.9** | 86.4 | 93.4 | <u>86.5</u> | 90.4 |
| UI-TARS-7B[*] | 93.1 | 89.8 | 98.2 | 87.9 | 94.7 | **91.0** | **93.0** |
| UI-TARS-1.5-7B[*] | 93.6 | 86.1 | <u>97.5</u> | 87.9 | 92.9 | **91.0** | 92.2 |
| **Phi-Ground-4B-7C** | 93.6 | 84.7 | 94.7 | 87.5 | 96.8 | 84.0 | 91.0 |
| **Phi-Ground-4B-16C** | 93.6 | 88.3 | 93.4 | **88.5** | 97.9 | 85.9 | 91.9 |
| **Phi-Ground-4B-16C-DPO** | 94.6 | 90.1 | <u>97.5</u> | 87.9 | <u>98.7</u> | 86.0 | 92.3 |
| **Phi-Ground-4B-29C** | 93.1 | **91.2** | 95.6 | 87.5 | 97.5 | 85.0 | 92.1 |
| **Phi-Ground-7B-7C** | 93.6 | 89.1 | 95.1 | 87.0 | 98.3 | 85.9 | 92.1 |
| **Phi-Ground-7B-16C** | 93.6 | <u>90.5</u> | 94.7 | <u>88.0</u> | 96.8 | 85.0 | 91.9 |
| **Phi-Ground-7B-16C-DPO** | **95.1** | 90.1 | 96.6 | **88.5** | **99.6** | 86.4 | **93.4** |
| **Phi-Ground-7B-29C** | <u>94.7</u> | 86.1 | 94.2 | **88.5** | 97.9 | 84.5 | 91.7 |

Table 13: ScreenSpot-Pro results.

| Model | Development | | | Creative | | | CAD | | | Scientific | | | Office | | | OS | | | AVG. | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Text | Icon | AVG. | Text | Icon | AVG. | Text | Icon | AVG. | Text | Icon | AVG. | Text | Icon | AVG. | Text | Icon | AVG. | Text | Icon | AVG. |
| *End-to-end model setting (Use short REs)* | | | | | | | | | | | | | | | | | | | | | |
| SeeClick [26] | 0.6 | 0.0 | 0.3 | 1.0 | 0.0 | 0.6 | 2.5 | 0.0 | 1.9 | 3.5 | 0.0 | 2.0 | 1.1 | 0.0 | 0.9 | 2.8 | 0.0 | 1.5 | 1.8 | 0.0 | 1.1 |
| OS-Atlas-4B [31] | 7.1 | 0.0 | 3.7 | 3.0 | 1.4 | 2.3 | 2.0 | 0.0 | 1.5 | 9.0 | 5.5 | 7.5 | 5.1 | 3.8 | 4.8 | 5.6 | 0.0 | 3.1 | 5.0 | 1.7 | 3.7 |
| Show-UI-2B | 16.9 | 1.4 | 9.4 | 9.1 | 0.0 | 5.3 | 2.5 | 0.0 | 1.9 | 13.2 | 7.3 | 10.6 | 15.3 | 7.5 | 13.5 | 10.3 | 2.2 | 6.6 | 10.8 | 2.6 | 7.7 |
| CogAgent-18B | 14.9 | 0.7 | 8.0 | 9.6 | 0.0 | 5.6 | 7.1 | 3.1 | 6.1 | 22.2 | 1.8 | 13.4 | 13.0 | 0.0 | 10.0 | 5.6 | 0.0 | 3.1 | 12.0 | 0.8 | 7.7 |
| Aria-UI | 16.2 | 0.0 | 8.4 | 23.7 | 2.1 | 14.7 | 7.6 | 1.6 | 6.1 | 27.1 | 6.4 | 18.1 | 20.3 | 1.9 | 16.1 | 4.7 | 0.0 | 2.6 | 17.1 | 2.0 | 11.3 |
| UGround-7B | 26.6 | 2.1 | 14.7 | 27.3 | 2.8 | 17.0 | 14.2 | 1.6 | 11.1 | 31.9 | 2.7 | 19.3 | 31.6 | 11.3 | 27.0 | 17.8 | 0.0 | 9.7 | 25.0 | 2.8 | 16.5 |
| Claude Computer Use [13] | 22.0 | 3.9 | 12.6 | 25.9 | 3.4 | 16.8 | 14.5 | 3.7 | 11.9 | 33.9 | 15.8 | 25.8 | 30.1 | 16.3 | 26.9 | 11.0 | 4.5 | 8.1 | 23.4 | 7.1 | 17.1 |
| OS-Atlas-7B | 33.1 | 1.4 | 17.7 | 28.8 | 2.8 | 17.9 | 12.2 | 4.7 | 10.3 | 37.5 | 7.3 | 24.4 | 33.9 | 5.7 | 27.4 | 27.1 | 4.5 | 16.8 | 28.1 | 4.0 | 18.9 |
| UI-TARS-2B | 47.4 | 4.1 | 26.4 | 42.9 | 6.3 | 27.6 | 17.8 | 4.7 | 14.6 | 56.9 | 17.3 | 39.8 | 50.3 | 17.0 | 42.6 | 21.5 | 5.6 | 14.3 | 39.6 | 8.4 | 27.7 |
| UGround-v1-7B | - | - | 35.5 | - | - | 27.8 | - | - | 13.5 | - | - | 38.8 | - | - | 48.8 | - | - | 26.1 | - | - | 31.1 |
| UI-TARS-7B | 58.4 | 12.4 | 36.1 | 50.0 | 9.1 | 32.8 | 20.8 | 9.4 | 18.0 | 63.9 | 31.8 | 50.0 | 63.3 | 20.8 | 53.5 | 30.8 | 16.9 | 24.5 | 47.8 | 16.2 | 35.7 |
| UI-TARS-1.5-7B* | 58.4 | 14.1 | 37.2 | 57.6 | 14.7 | 39.6 | 39.1 | 15.6 | 33.3 | 68.3 | 23.4 | 49.0 | 73.6 | 34.7 | 65.0 | 49.5 | 17.2 | 35.1 | 57.5 | 18.2 | 42.6 |
| **Phi-Ground-4B-7C** | 32.5 | 5.6 | 19.5 | 26.8 | 7.0 | 18.5 | 10.2 | 4.7 | 8.8 | 39.9 | 17.3 | 30.0 | 46.6 | 34.0 | 43.8 | 27.1 | 16.9 | 22.4 | 29.8 | 12.0 | 23.0 |
| **Phi-Ground-4B-16C** | 59.1 | 18.8 | 39.6 | 48.0 | 9.1 | 31.7 | 15.7 | 9.4 | 14.2 | 53.8 | 25.5 | 41.5 | 69.0 | 42.0 | 62.9 | 45.8 | 21.3 | 34.7 | 47.5 | 19.0 | 36.7 |
| **Phi-Ground-4B-16C-DPO** | 64.3 | 15.3 | 40.6 | 46.5 | 12.6 | 32.3 | 20.3 | 12.5 | 18.4 | 54.5 | 26.4 | 42.3 | 69.0 | 42.0 | 62.9 | 44.9 | 24.7 | 35.7 | 49.0 | 20.0 | 38.0 |
| **Phi-Ground-4B-29C** | 64.9 | 15.3 | 40.9 | 48.5 | 7.7 | 31.4 | 21.8 | 9.4 | 18.8 | 58.7 | 23.6 | 43.5 | 71.8 | 36.0 | 63.8 | 47.7 | 18.0 | 34.2 | 51.3 | 16.5 | 38.0 |
| **Phi-Ground-7B-7C** | 57.8 | 12.5 | 35.9 | 39.9 | 9.8 | 27.3 | 17.8 | 6.3 | 14.9 | 53.2 | 24.6 | 40.7 | 63.2 | 22.0 | 54.0 | 42.1 | 19.1 | 31.6 | 44.6 | 15.2 | 33.4 |
| **Phi-Ground-7B-16C** | 68.2 | 16.7 | 43.3 | 53.5 | 11.2 | 35.8 | 26.4 | 17.2 | 24.1 | 58.7 | 25.5 | 44.3 | 76.4 | 42.0 | 68.8 | 53.3 | 22.5 | 39.3 | 55.2 | 20.0 | 41.8 |
| **Phi-Ground-7B-16C-DPO** | 70.8 | 16.7 | 44.6 | 56.6 | 13.3 | 38.4 | 26.9 | 17.2 | 24.5 | 58.0 | 29.1 | 45.5 | 76.4 | 44.0 | 69.2 | 55.1 | 25.8 | 41.8 | 56.4 | 21.8 | 43.2 |
| **Phi-Ground-7B-29C** | 66.9 | 20.1 | 44.3 | 57.1 | 9.8 | 37.2 | 26.9 | 10.9 | 23.0 | 60.5 | 25.5 | 47.8 | 79.3 | 40.0 | 70.5 | 57.9 | 38.0 | 19.8 | 57.8 | 19.0 | 43.0 |
| *Agent setting (Use long REs) with GPT-4O as planner* | | | | | | | | | | | | | | | | | | | | | |
| SeeClick* | 1.3 | 1.4 | 1.3 | 1.0 | 0.0 | 0.6 | 1.0 | 0.0 | 0.8 | 2.8 | 0.0 | 1.6 | 0.6 | 0.0 | 0.4 | 0.0 | 0.0 | 0.0 | 1.1 | 0.3 | 0.8 |
| UGround-7B* | 33.1 | 2.1 | 18.1 | 33.3 | 4.9 | 21.4 | 14.7 | 1.6 | 11.5 | 40.6 | 5.5 | 25.3 | 37.9 | 8.0 | 31.2 | 24.3 | 2.2 | 14.3 | 30.4 | 3.8 | 20.3 |
| UGround-v1-7B* | 46.8 | 2.8 | 25.5 | 37.4 | 4.9 | 23.8 | 16.8 | 3.1 | 13.4 | 53.8 | 15.5 | 37.2 | 53.4 | 18.0 | 45.5 | 34.6 | 5.6 | 21.4 | 39.7 | 7.3 | 27.3 |
| OS-Atlas-4B* | 2.6 | 0.0 | 1.3 | 2.5 | 0.0 | 1.5 | 1.5 | 0.0 | 1.1 | 3.5 | 0.0 | 2.0 | 2.3 | 2.0 | 2.2 | 2.8 | 0.0 | 1.5 | 2.5 | 0.2 | 1.6 |
| OS-Atlas-7B* | 34.5 | 0.0 | 17.5 | 29.8 | 3.5 | 18.8 | 7.2 | 3.3 | 6.3 | 36.6 | 7.4 | 24.0 | 35.3 | 10.4 | 29.8 | 19.6 | 6.9 | 13.9 | 26.7 | 4.4 | 18.2 |
| UI-TARS-2B* | 48.1 | 6.2 | 27.9 | 42.4 | 9.1 | 28.4 | 13.2 | 3.1 | 10.7 | 65.0 | 14.5 | 43.1 | 53.4 | 18.0 | 45.5 | 29.0 | 11.2 | 20.9 | 41.2 | 9.8 | 29.2 |
| UI-TARS-7B* | 56.5 | 10.4 | 34.2 | 49.0 | 14.0 | 34.3 | 18.3 | 4.7 | 14.9 | 69.2 | 24.5 | 49.8 | 61.5 | 20.0 | 52.2 | 33.6 | 13.5 | 24.5 | 47.5 | 14.5 | 34.9 |
| UI-TARS-1.5-7B* | 57.1 | 13.2 | 35.9 | 49.5 | 14.7 | 34.9 | 28.4 | 12.5 | 24.5 | 67.8 | 23.6 | 48.6 | 67.8 | 42.0 | 62.1 | 40.2 | 18.0 | 30.1 | 51.4 | 18.5 | 38.8 |
| **Phi-Ground-4B-7C** | 42.9 | 6.3 | 25.2 | 33.8 | 13.3 | 25.2 | 15.2 | 1.6 | 11.9 | 55.9 | 16.4 | 38.7 | 48.3 | 36.0 | 15.5 | 33.6 | 12.4 | 24.0 | 37.3 | 12.7 | 27.9 |
| **Phi-Ground-4B-16C** | 49.4 | 16.0 | 33.2 | 49.0 | 16.8 | 35.5 | 17.8 | 15.6 | 17.2 | 59.4 | 20.9 | 42.7 | 62.1 | 50.0 | 58.9 | 46.8 | 21.3 | 36.2 | 46.6 | 20.7 | 36.7 |
| **Phi-Ground-4B-16C-DPO** | 48.1 | 16.0 | 32.6 | 52.0 | 18.9 | 38.1 | 18.8 | 14.1 | 17.6 | 61.5 | 25.5 | 45.8 | 63.2 | 44.0 | 58.9 | 49.5 | 24.7 | 38.3 | 47.8 | 21.8 | 37.9 |
| **Phi-Ground-4B-29C** | 54.5 | 17.4 | 36.6 | 50.0 | 14.0 | 34.3 | 20.8 | 7.8 | 17.6 | 67.8 | 20.9 | 47.4 | 67.8 | 34.0 | 60.3 | 44.9 | 15.7 | 31.6 | 50.1 | 17.3 | 37.8 |
| **Phi-Ground-7B-7C** | 74.2 | 20.3 | 48.0 | 53.3 | 19.0 | 38.9 | 25.0 | 14.1 | 22.3 | 68.3 | 26.7 | 50.6 | 75.2 | 33.3 | 66.0 | 58.8 | 20.5 | 41.1 | 57.3 | 21.5 | 43.6 |
| **Phi-Ground-7B-16C** | 53.2 | 16.0 | 35.2 | 52.0 | 18.2 | 37.8 | 21.8 | 12.5 | 19.5 | 68.5 | 22.7 | 48.6 | 63.2 | 42.0 | 58.4 | 42.1 | 14.6 | 29.6 | 49.4 | 17.8 | 38.0 |
| **Phi-Ground-7B-16C-DPO** | 55.8 | 16.7 | 36.9 | 54.0 | 20.3 | 39.9 | 25.4 | 17.2 | 23.4 | 65.7 | 25.5 | 48.2 | 65.5 | 46.0 | 61.2 | 43.0 | 18.0 | 31.6 | 51.1 | 21.8 | 39.9 |
| **Phi-Ground-7B-29C** | 55.8 | 18.8 | 37.9 | 56.6 | 15.4 | 39.3 | 24.4 | 14.1 | 21.8 | 68.5 | 25.5 | 49.8 | 61.5 | 32.0 | 54.9 | 48.6 | 16.9 | 34.2 | 51.7 | 19.5 | 39.4 |
| *Agent setting (Use long REs) with O4-mini as planner* | | | | | | | | | | | | | | | | | | | | | |
| SeeClick* | 3.3 | 0.0 | 1.7 | 0.5 | 0.0 | 0.3 | 1.5 | 0.0 | 1.2 | 2.1 | 1.0 | 1.6 | 3.7 | 0.0 | 2.9 | 0.0 | 0.0 | 0.0 | 1.9 | 0.2 | 1.2 |
| UGround-7B* | 41.1 | 4.9 | 23.5 | 34.4 | 5.6 | 22.3 | 14.8 | 3.1 | 11.9 | 40.8 | 9.5 | 27.5 | 47.2 | 11.1 | 39.3 | 35.3 | 3.4 | 20.5 | 34.6 | 6.0 | 23.7 |
| UGround-v1-7B* | 62.9 | 4.9 | 34.7 | 42.1 | 7.0 | 27.3 | 21.4 | 6.2 | 17.7 | 51.4 | 19.0 | 37.7 | 60.9 | 20.0 | 51.9 | 47.1 | 11.4 | 30.5 | 46.3 | 10.2 | 32.5 |
| OS-Atlas-4B* | 3.3 | 0.7 | 2.0 | 2.1 | 0.0 | 1.2 | 2.1 | 0.0 | 1.5 | 4.2 | 0.0 | 2.4 | 4.4 | 2.2 | 3.9 | 2.9 | 0.0 | 1.6 | 3.1 | 0.3 | 2.0 |
| OS-Atlas-7B* | 41.4 | 1.4 | 21.6 | 30.4 | 5.6 | 19.9 | 10.8 | 8.1 | 10.1 | 38.0 | 13.5 | 27.6 | 37.3 | 13.3 | 32.0 | 25.7 | 6.8 | 16.9 | 29.8 | 7.0 | 21.1 |
| UI-TARS-2B* | 64.2 | 6.3 | 36.1 | 51.8 | 10.6 | 34.4 | 18.9 | 9.4 | 16.5 | 63.4 | 23.8 | 46.6 | 65.8 | 20.0 | 55.8 | 43.1 | 12.5 | 28.9 | 50.2 | 12.8 | 35.9 |
| UI-TARS-7B* | 66.2 | 11.9 | 39.8 | 55.9 | 14.1 | 38.3 | 25.0 | 10.9 | 21.5 | 72.5 | 28.6 | 53.8 | 71.4 | 22.2 | 60.7 | 51.0 | 12.5 | 33.2 | 55.8 | 16.2 | 40.6 |
| UI-TARS-1.5-7B* | 67.5 | 23.8 | 46.3 | 59.5 | 23.9 | 44.5 | 43.4 | 20.3 | 37.7 | 70.4 | 29.5 | 53.0 | 83.2 | 42.2 | 74.3 | 58.8 | 22.7 | 42.1 | 63.0 | 25.7 | 48.8 |
| **Phi-Ground-4B-7C** | 69.5 | 14.7 | 42.9 | 44.6 | 20.4 | 34.4 | 23.0 | 4.7 | 18.5 | 69.7 | 23.8 | 50.2 | 72.0 | 51.1 | 67.5 | 52.0 | 22.7 | 38.4 | 53.3 | 20.6 | 40.8 |
| **Phi-Ground-4B-16C** | 74.2 | 29.4 | 52.4 | 59.5 | 31.0 | 47.5 | 29.6 | 18.8 | 26.9 | 69.0 | 27.6 | 51.4 | 82.6 | 48.9 | 75.2 | 67.6 | 33.0 | 51.6 | 61.9 | 30.3 | 49.8 |
| **Phi-Ground-4B-16C-DPO** | 76.2 | 33.6 | 55.4 | 62.1 | 31.9 | 49.4 | 32.1 | 15.6 | 28.1 | 69.0 | 31.4 | 53.0 | 82.6 | 53.3 | 76.2 | 66.7 | 36.4 | 52.6 | 63.1 | 32.8 | 51.5 |
| **Phi-Ground-4B-29C** | 76.2 | 22.4 | 50.0 | 63.1 | 27.5 | 48.1 | 29.6 | 14.1 | 25.8 | 76.8 | 28.6 | 56.3 | 83.2 | 40.0 | 73.8 | 63.7 | 20.5 | 43.7 | 63.8 | 24.9 | 48.9 |
| **Phi-Ground-7B-7C** | 52.0 | 14.6 | 33.9 | 41.4 | 11.9 | 29.0 | 18.8 | 7.8 | 16.1 | 62.2 | 20.0 | 43.9 | 58.1 | 28.0 | 51.3 | 35.5 | 14.6 | 26.0 | 43.9 | 15.3 | 33.0 |
| **Phi-Ground-7B-16C** | 76.8 | 32.2 | 55.1 | 66.2 | 35.9 | 53.4 | 34.7 | 17.2 | 30.4 | 73.9 | 35.2 | 57.5 | 83.2 | 51.1 | 76.2 | 69.6 | 31.8 | 52.1 | 65.8 | 33.4 | 53.4 |
| **Phi-Ground-7B-16C-DPO** | 80.8 | 32.9 | 57.5 | 66.7 | 37.3 | 54.3 | 35.7 | 18.8 | 31.5 | 77.5 | 35.2 | 59.5 | 83.9 | 57.8 | 78.2 | 71.6 | 31.8 | 53.2 | 67.6 | 34.6 | 55.0 |
| **Phi-Ground-7B-29C** | 78.1 | 29.4 | 54.4 | 67.2 | 31.0 | 51.9 | 40.3 | 18.8 | 35.0 | 77.5 | 35.2 | 59.5 | 83.2 | 51.1 | 76.2 | 72.5 | 30.7 | 53.2 | 68.2 | 31.5 | 54.2 |

| Model | Basic | | | | | | | Functional | | | | | | | Spatial | | | | | | | Final Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Ed (215) | Br (56) | De (376) | Pr (605) | Cr (438) | En (82) | Overall (1772) | Ed (215) | Br (56) | De (376) | Pr (605) | Cr (438) | En (82) | Overall (1772) | Ed (212) | Br (31) | De (338) | Pr (740) | Cr (586) | En (28) | Overall (1935) | |
| *Closed-Source VLMs* | | | | | | | | | | | | | | | | | | | | | | |
| GPT-4o | 2.23 | 0.00 | 1.86 | 1.16 | 1.14 | 4.88 | 1.58 | 1.40 | 0.00 | 3.19 | 0.83 | 0.91 | 3.66 | 1.52 | 0.94 | 0.00 | 1.48 | 1.22 | 0.51 | 3.57 | 1.03 | 1.38 |
| Gemini-1.5-pro | 0.47 | 0.00 | 1.60 | 0.83 | 0.46 | 0.00 | 0.79 | 0.00 | 1.79 | 0.27 | 0.17 | 0.46 | 0.00 | 0.28 | 0.94 | 0.00 | 0.89 | 0.54 | 0.34 | 0.00 | 0.57 | 0.55 |
| Gemini-Flash-2.0 | 0.00 | 0.00 | 0.27 | 0.66 | 0.68 | 0.00 | 0.45 | 0.47 | 1.79 | 0.00 | 0.66 | 0.23 | 0.00 | 0.40 | 0.00 | 0.00 | 0.00 | 0.14 | 0.00 | 0.00 | 0.05 | 0.30 |
| Claude-3.5-Sonnet | 3.26 | 16.1 | 5.32 | 6.94 | 1.83 | 4.88 | 5.08 | 5.12 | 19.6 | 4.79 | 5.95 | 2.51 | 4.88 | 5.19 | 2.83 | 9.68 | 5.03 | 2.43 | 2.56 | 7.14 | 3.15 | 4.47 |
| Claude-3.7-Sonnet | 6.51 | 12.5 | 7.98 | 11.24 | 9.13 | 11.0 | 9.48 | 5.12 | 7.14 | 8.24 | 9.92 | 6.16 | 4.88 | 7.73 | 6.60 | 9.68 | 7.69 | 7.43 | 7.85 | 10.7 | 7.60 | 8.27 |
| *Open-Source VLMs* | | | | | | | | | | | | | | | | | | | | | | |
| Qwen-2.5VL-7B | 0.47 | 0.00 | 1.33 | 1.65 | 0.68 | 1.22 | 1.24 | 0.47 | 0.00 | 0.80 | 1.16 | 0.46 | 1.22 | 0.79 | 0.47 | 0.00 | 1.48 | 0.00 | 0.51 | 0.00 | 0.51 | 0.85 |
| InternVL2-8B | 0.00 | 0.00 | 0.00 | 0.02 | 0.00 | 0.14 | 0.11 | 0.00 | 0.00 | 0.27 | 0.00 | 0.00 | 1.22 | 0.11 | 0.00 | 0.00 | 0.00 | 0.14 | 0.00 | 0.00 | 0.05 | 0.09 |
| InternVL2.5-8B | 0.93 | 8.93 | 3.46 | 2.31 | 1.37 | 4.88 | 2.48 | 1.40 | 7.14 | 3.72 | 2.81 | 1.60 | 6.10 | 2.82 | 0.94 | 3.23 | 1.78 | 0.68 | 0.68 | 3.57 | 0.98 | 2.09 |
| Qwen-2VL-7B | 2.79 | 7.14 | 3.72 | 3.97 | 0.68 | 12.2 | 3.44 | 2.79 | 12.5 | 3.19 | 3.97 | 0.68 | 6.10 | 3.22 | 0.47 | 3.23 | 2.37 | 1.08 | 0.51 | 3.57 | 1.45 | 2.70 |
| MiniCPM-V-8B | 4.19 | 21.4 | 7.71 | 7.44 | 3.65 | 18.3 | 7.11 | 4.19 | 19.6 | 6.38 | 4.63 | 2.97 | 11.0 | 5.30 | 0.47 | 3.23 | 1.78 | 0.27 | 0.17 | 3.57 | 1.45 | 4.34 |
| *Open-Source GUI Agents (≤ 10B)* | | | | | | | | | | | | | | | | | | | | | | |
| ShowUI-2B | 5.12 | 16.1 | 9.84 | 9.09 | 3.42 | 19.5 | 8.07 | 5.12 | 12.5 | 9.31 | 8.60 | 4.11 | 15.9 | 7.67 | 0.94 | 9.68 | 2.96 | 2.70 | 0.68 | 3.57 | 2.07 | 5.94 |
| AriaUI-25.3B | 10.7 | 23.2 | 13.0 | 12.9 | 8.22 | 20.7 | 12.2 | 12.6 | 19.6 | 15.4 | 14.6 | 10.5 | 22.0 | 14.0 | 3.77 | 9.68 | 4.44 | 4.86 | 2.22 | 7.14 | 3.98 | 10.1 |
| UGround-v1-7B | 11.6 | 35.7 | 19.7 | 15.0 | 11.0 | 18.3 | 15.4 | 15.4 | 33.9 | 22.3 | 16.5 | 11.6 | 19.5 | 17.1 | 4.25 | 6.45 | 9.76 | 6.35 | 4.44 | 14.3 | 6.25 | 12.9 |
| OSAtlas-7B | 10.7 | 23.2 | 13.3 | 12.6 | 8.22 | 22.0 | 12.2 | 11.6 | 16.1 | 11.4 | 12.7 | 7.53 | 13.4 | 11.2 | 3.77 | 6.45 | 5.62 | 3.65 | 2.22 | 7.14 | 3.67 | 9.02 |
| UGround-7B | 10.2 | 23.2 | 14.9 | 10.6 | 7.53 | 19.5 | 11.5 | 12.1 | 25.0 | 15.2 | 11.2 | 7.99 | 20.7 | 12.2 | 2.36 | 0.00 | 4.14 | 2.70 | 2.22 | 7.14 | 2.79 | 8.83 |
| Aguvis-7B | 16.7 | 37.5 | 22.3 | 16.2 | 12.6 | 26.8 | 17.8 | 17.2 | 35.7 | 21.5 | 18.0 | 13.0 | 24.4 | 18.3 | 5.19 | 9.68 | 6.51 | 4.05 | 4.78 | 14.3 | 5.06 | 13.7 |
| UI-TARS-7B | 15.4 | 41.1 | 21.8 | 21.2 | 13.2 | 39.0 | 20.1 | 20.5 | 41.1 | 25.5 | 26.5 | 16.0 | 45.1 | 24.3 | 6.60 | 12.9 | 11.0 | 9.19 | 5.80 | 17.9 | 8.37 | 17.6 |
| UI-TARS-1.5-7B * | 22.1 | 59.6 | 29.5 | 31.8 | 19.8 | 50.0 | 28.8 | 24.3 | 56.9 | 26.2 | 31.2 | 19.1 | 40.7 | 27.5 | 8.96 | 20.0 | 10.9 | 11.6 | 8.70 | 29.6 | 10.7 | 22.3 |
| CogAgent-9B | 11.2 | 14.3 | 12.5 | 13.7 | 8.68 | 15.9 | 12.0 | 11.6 | 14.3 | 11.4 | 14.7 | 8.22 | 18.3 | 12.2 | 3.30 | 0.00 | 1.18 | 4.05 | 1.37 | 7.14 | 2.63 | 8.94 |
| SeeClick-9.6B | 7.44 | 23.2 | 13.0 | 8.43 | 5.48 | 17.1 | 9.42 | 4.65 | 7.14 | 5.32 | 3.97 | 4.34 | 7.32 | 4.68 | 0.47 | 6.45 | 3.25 | 1.22 | 2.73 | 3.57 | 2.07 | 5.39 |
| *Our Phi-Ground Models* | | | | | | | | | | | | | | | | | | | | | | |
| **Phi-Ground-4B-7C** | 26.5 | 42.9 | 27.2 | 31.6 | 18.7 | 38.3 | 27.5 | 27.4 | 46.4 | 26.9 | 34.9 | 21.0 | 35.8 | 29.2 | 6.13 | 6.45 | 5.92 | 4.60 | 2.90 | 18.5 | 4.71 | 20.5 |
| **Phi-Ground-4B-16C** | 27.0 | 50.0 | 30.9 | 37.2 | 24.4 | 40.7 | 32.0 | 28.4 | 48.2 | 30.2 | 39.3 | 21.2 | 43.2 | 32.5 | 5.66 | 3.23 | 7.82 | 3.58 | 3.70 | | 5.28 | 23.3 |
| **Phi-Ground-4B-16C-DPO** | 29.8 | 48.2 | 30.7 | 39.7 | 24.9 | 45.7 | 33.4 | 32.1 | 51.8 | 32.4 | 41.2 | 23.1 | 45.7 | 34.3 | 6.13 | 0.00 | 9.17 | 5.68 | 4.27 | 3.70 | 5.79 | 24.5 |
| **Phi-Ground-4B-29C** | 29.3 | 53.6 | 32.3 | 34.5 | 22.4 | 48.1 | 31.6 | 27.0 | 46.4 | 31.6 | 38.5 | 20.1 | 37.0 | 31.3 | 7.08 | 3.23 | 4.14 | 3.65 | 2.56 | 7.41 | 3.83 | 22.2 |
| **Phi-Ground-7B-7C** | 28.4 | 53.6 | 32.8 | 34.7 | 23.7 | 45.7 | 31.9 | 25.6 | 53.6 | 33.0 | 37.2 | 21.0 | 39.5 | 31.5 | 6.11 | 9.74 | 5.31 | 5.80 | 5.60 | 14.8 | 5.90 | 23.1 |
| **Phi-Ground-7B-16C** | 35.3 | 53.6 | 38.9 | 38.8 | 24.0 | 50.5 | 35.8 | 32.6 | 62.5 | 36.7 | 38.7 | 23.7 | 43.2 | 34.8 | 6.13 | 3.23 | 7.40 | 6.08 | 6.32 | 3.70 | 6.31 | 25.6 |
| **Phi-Ground-7B-16C-DPO** | 35.3 | 57.1 | 37.1 | 40.3 | 26.0 | 56.8 | 36.8 | 34.4 | 62.5 | 39.1 | 40.8 | 25.8 | 50.6 | 37.1 | 7.51 | 6.50 | 8.33 | 6.80 | 8.41 | 3.70 | 7.60 | 27.2 |
| **Phi-Ground-7B-29C** | 33.5 | 46.4 | 36.3 | 34.0 | 23.3 | 51.9 | 33.0 | 28.8 | 51.8 | 36.2 | 34.7 | 25.8 | 43.2 | 33.0 | 8.02 | 6.45 | 7.69 | 4.86 | 4.78 | 0.00 | 5.63 | 23.9 |
| *GPT-4O as planner* | | | | | | | | | | | | | | | | | | | | | | |
| UGround-v1-7B* | 21.9 | 44.6 | 27.2 | 21.8 | 17.6 | 25.9 | 22.8 | 26.0 | 33.9 | 28.5 | 22.0 | 16.9 | 27.2 | 23.2 | 7.08 | 9.68 | 13.0 | 7.31 | 6.48 | 14.8 | 8.17 | 18.1 |
| OSAtlas-7B* | 11.7 | 30.4 | 19.1 | 14.1 | 8.03 | 22.5 | 14.3 | 13.1 | 30.4 | 18.1 | 13.5 | 9.40 | 21.0 | 14.3 | 4.74 | 9.68 | 8.04 | 3.52 | 3.95 | 3.70 | 4.68 | 11.1 |
| UGround-7B* | 17.2 | 25.0 | 22.4 | 17.9 | 12.6 | 24.7 | 18.0 | 19.5 | 26.8 | 22.1 | 18.5 | 12.6 | 24.7 | 18.5 | 6.13 | 9.68 | 10.4 | 5.28 | 5.63 | 3.70 | 6.41 | 14.3 |
| UI-TARS-7B* | 23.3 | 44.6 | 28.8 | 27.6 | 19.4 | 46.9 | 26.7 | 24.7 | 42.9 | 29.0 | 28.1 | 16.9 | 37.0 | 26.0 | 12.3 | 9.68 | 14.2 | 10.8 | 5.98 | 22.2 | 10.2 | 21.0 |
| UI-TARS-1.5-7B* | 24.9 | 47.3 | 32.3 | 29.3 | 22.1 | 50.6 | 29.1 | 26.5 | 47.3 | 28.2 | 31.0 | 20.5 | 39.5 | 28.2 | 9.43 | 9.68 | 16.0 | 9.47 | 9.04 | 18.5 | 10.6 | 22.6 |
| SeeClick-9.6B* | 3.72 | 8.93 | 4.01 | 3.80 | 4.34 | 3.70 | 4.12 | 3.72 | 5.36 | 4.52 | 3.47 | 3.21 | 3.70 | 3.73 | 0.47 | 3.23 | 3.55 | 1.49 | 1.71 | 3.70 | 1.86 | 3.24 |
| **Phi-Ground-4B-7C** | 26.0 | 50.0 | 32.0 | 32.9 | 23.7 | 48.1 | 30.8 | 28.4 | 42.9 | 32.4 | | 21.5 | 42.0 | 29.6 | 9.91 | 9.68 | 16.3 | 10.1 | 8.70 | 14.8 | 10.9 | 23.8 |
| **Phi-Ground-4B-16C** | 27.0 | 46.4 | 33.3 | 35.7 | 26.5 | 53.1 | 33.0 | 29.3 | 41.1 | 32.7 | 37.9 | 23.7 | 48.2 | 32.8 | 11.3 | 16.1 | 15.7 | 9.91 | 8.01 | 11.1 | 10.6 | 25.5 |
| **Phi-Ground-4B-16C-DPO** | 26.2 | 51.8 | 33.6 | 38.1 | 26.0 | 56.8 | 34.0 | 29.3 | 41.1 | 34.0 | 39.1 | 23.8 | 45.0 | 33.4 | 10.4 | 12.9 | 18.3 | 12.6 | 7.40 | 18.5 | 11.9 | 26.4 |
| **Phi-Ground-4B-29C** | 28.4 | 48.2 | 33.3 | 36.5 | 26.5 | 48.1 | 33.3 | 29.8 | 44.6 | 33.5 | 36.0 | 24.9 | 48.1 | 32.8 | 11.8 | 12.9 | 15.7 | 11.0 | 9.04 | 3.70 | 11.2 | 25.8 |
| **Phi-Ground-7B-7C** | 30.7 | 51.8 | 33.1 | 32.4 | 27.4 | 51.9 | 32.6 | 29.3 | 50.0 | 32.4 | 33.2 | 24.0 | 42.0 | 31.2 | 11.8 | 16.1 | 11.2 | 9.30 | 9.21 | 14.8 | 10.1 | 24.6 |
| **Phi-Ground-7B-16C** | 31.6 | 51.8 | 35.2 | | 28.1 | 54.3 | 34.6 | 29.8 | 51.8 | 37.5 | 36.2 | 24.2 | 49.4 | 33.8 | 12.3 | 16.1 | 14.8 | 10.5 | 9.74 | 11.1 | 11.3 | 26.6 |
| **Phi-Ground-7B-16C-DPO** | 30.7 | 50.0 | 38.7 | 38.0 | 27.9 | 61.7 | 36.2 | 32.1 | 50.0 | 39.1 | 39.0 | 25.1 | 50.6 | 35.6 | 12.7 | 19.4 | 15.1 | 11.6 | 9.40 | 14.8 | 11.9 | 27.9 |
| **Phi-Ground-7B-29C** | 31.6 | 48.2 | 36.7 | 33.4 | 28.8 | 60.5 | 34.0 | 29.3 | 46.4 | 36.4 | 35.7 | 26.7 | 49.4 | 33.8 | 11.3 | 19.4 | 13.6 | 10.9 | 9.91 | 14.8 | 11.3 | 26.4 |
| *O4-mini as planner* | | | | | | | | | | | | | | | | | | | | | | |
| UGround-v1-7B* | 31.3 | 56.4 | 32.8 | 29.4 | 24.0 | 43.8 | 30.5 | 32.5 | 44.6 | 32.5 | 27.2 | 23.2 | 42.5 | 29.3 | 11.3 | 16.1 | 19.8 | 12.6 | 11.8 | 37.0 | 13.9 | 24.6 |
| OSAtlas-7B* | 16.4 | 32.7 | 21.9 | 15.8 | 11.0 | 30.0 | 17.2 | 17.9 | 30.4 | 20.5 | 13.9 | 11.6 | 33.3 | 16.6 | 6.67 | 12.9 | 13.4 | 5.27 | 5.87 | 29.6 | 7.48 | 13.8 |
| UGround-7B* | 23.4 | 41.8 | 29.0 | 23.3 | 19.5 | 45.0 | 25.1 | 25.5 | 33.9 | 26.8 | 19.6 | 17.4 | 42.5 | 22.8 | 10.4 | 19.4 | 15.1 | 10.3 | 9.76 | 22.2 | 11.3 | 19.7 |
| UI-TARS-7B* | 30.4 | 54.5 | 36.3 | 32.5 | 26.3 | 53.8 | 33.2 | 35.8 | 50.0 | 35.2 | 34.6 | 23.2 | 51.2 | 33.4 | 16.0 | 25.8 | 22.5 | 16.9 | 11.6 | 44.4 | 16.7 | 27.8 |
| UI-TARS-1.5-7B * | 30.8 | 61.8 | 36.6 | 35.4 | 27.4 | 60.0 | 35.1 | 36.8 | 55.4 | 35.2 | 36.5 | 24.9 | 60.0 | 35.1 | 17.5 | 25.8 | 23.7 | 16.2 | 15.2 | 44.4 | 17.9 | 29.4 |
| SeeClick-9.6B* | 5.61 | 10.9 | 5.46 | 5.15 | 3.72 | 7.50 | 5.21 | 4.25 | 7.14 | 6.56 | 4.57 | 3.52 | 10.0 | 5.03 | 1.42 | 9.68 | 3.85 | 2.84 | 3.08 | 11.1 | 3.16 | 4.47 |
| **Phi-Ground-4B-7C** | 35.5 | 60.0 | 36.3 | 40.2 | 27.7 | 60.0 | 37.3 | 39.2 | 58.9 | 36.3 | 39.6 | 26.3 | 60.0 | 37.1 | 16.5 | 22.6 | 25.4 | 19.2 | 13.4 | 44.4 | 18.6 | 31.0 |
| **Phi-Ground-4B-16C** | 37.9 | 63.6 | 40.2 | 46.7 | 32.8 | 63.8 | 42.1 | 44.8 | 60.7 | 38.0 | 43.5 | 30.8 | 60.0 | 40.7 | 17.9 | 29.0 | 24.3 | 20.8 | 15.1 | 29.6 | 19.6 | 34.1 |
| **Phi-Ground-4B-16C-DPO** | 37.4 | 63.6 | 41.8 | 48.3 | 35.5 | 67.5 | 43.8 | 67.6 | 62.5 | 39.1 | 44.6 | 31.9 | 63.8 | 42.1 | 19.0 | 29.0 | 26.9 | 22.3 | 15.1 | 40.7 | 21.0 | 35.6 |
| **Phi-Ground-4B-29C** | 36.4 | 61.8 | 40.7 | 45.7 | 32.3 | 67.5 | 41.7 | 39.2 | 60.7 | 40.2 | 41.6 | 29.1 | 63.8 | 39.6 | 17.5 | 32.3 | 26.3 | 20.0 | 14.9 | 37.0 | 19.7 | 33.7 |
| **Phi-Ground-7B-7C** | 38.3 | 65.5 | 41.3 | 40.4 | 32.6 | 63.7 | 40.2 | 46.5 | 52.4 | 40.2 | 37.6 | 37.6 | 45.0 | 39.7 | 17.9 | 29.0 | 27.2 | 19.5 | 13.9 | 44.4 | 19.5 | 33.1 |
| **Phi-Ground-7B-16C** | 36.9 | 65.5 | 44.0 | 46.0 | 35.1 | 66.3 | 43.3 | 42.3 | 45.5 | 42.6 | 42.0 | 43.4 | 45.0 | 42.8 | 17.0 | 29.0 | 27.5 | 21.0 | 15.6 | 48.2 | 20.6 | 35.6 |
| **Phi-Ground-7B-16C-DPO** | 38.8 | 67.3 | 44.0 | 47.0 | 35.3 | 70.9 | 44.2 | 45.1 | 49.1 | 44.0 | 43.7 | 43.0 | 42.5 | 43.8 | 17.0 | 29.0 | 26.6 | 21.5 | 15.2 | 48.1 | 20.5 | 36.2 |
| **Phi-Ground-7B-29C** | 38.8 | 63.6 | 42.9 | 43.7 | 35.3 | 70.0 | 42.7 | 40.4 | 49.1 | 43.7 | 39.3 | 42.3 | 40.0 | 41.4 | 13.7 | 29.0 | 26.6 | 20.5 | 14.6 | 40.7 | 19.5 | 34.5 |

Table 14: Results of UI-Vision [24]. The final column shows the overall average. Abbreviated category labels: Ed (Education), Br (Browsers), De (Development), Pr (Productivity), Cr (Creativity), En (Entertainment). The best model within each size category is highlighted in **bold**, and the runner-up is underlined. We tested Agent setting results of UI-Vision of several open-source GUI models (≤ 10B parameters).

Table 15: Showdown-click-dev results. †: For the latency of the models we tested, we report the inference speed of the models accelerated using the vllm Python library if supported, otherwise we report the latency using huggingface transformers, marked with 'hf'. This may be faster than the results provided by the benchmark itself, but the comparison between the models we tested remains fair. For the settings with GPT-4O and O4-mini as planners, we directly added 2.5 seconds (aligned with the original benchmark) and 8 seconds (our tested average level, which may be highly dependent on the endpoint) to the original model latency, respectively. *: Results from the original GitHub repository.

| Model | Accuracy(%) ↑ | Latency†(ms) ↓ | Model | Accuracy(%) ↑ | Latency†(ms) ↓ |
|---|---|---|---|---|---|
| *End-to-end model setting (Use short REs)* | | | | | |
| GPT-4O* | 5.21 | 2500 | UI-TARS-1.5-7B | **67.15** | 445 |
| Qwen2.5-VL-73B-Instruct* | 24.78 | 3790 | UGround-7B | 46.50 | 1871 (hf) |
| Gemini 2.0 Flash* | 33.39 | 3069 | UGround-v1-7B | 57.81 | 209 |
| UI-TARS-72B-SFT* | 54.40 | 1977 | **Phi-Ground-4B-7C** | 54.40 | 122 |
| Claude 3.7 Sonnet (Computer Use)* | 53.68 | 9656 | **Phi-Ground-4B-16C** | 59.96 | 212 |
| Molmo-72B-0924* | 54.76 | 6599 | **Phi-Ground-4B-16C-DPO** | 58.17 | 212 |
| Operator (OpenAI CUA)* | 64.27 | 6385 | **Phi-Ground-4B-29C** | 55.11 | 401 |
| seeclick | 24.60 | 847 (hf) | **Phi-Ground-7B-7C** | 57.45 | 168 |
| OS-ATLAS-4B | 15.80 | 1288 (hf) | **Phi-Ground-7B-16C** | 61.04 | 313 |
| OS-ATLAS-7B | 41.11 | 1788 (hf) | **Phi-Ground-7B-16C-DPO** | 62.48 | 313 |
| UI-TARS-2B | 59.78 | 186 | **Phi-Ground-7B-29C** | 61.22 | 603 |
| UI-TARS-7B | *66.07* | 237 | | | |
| *Agent setting (Use long REs) with GPT-4O as planner* | | | | | |
| seeclick | 15.62 | 3347 (hf) | **Phi-Ground-4B-7C** | 60.68 | 2622 |
| OS-ATLAS-4B | 13.46 | 3788 (hf) | **Phi-Ground-4B-16C** | 62.84 | 2712 |
| OS-ATLAS-7B | 40.22 | 4288 (hf) | **Phi-Ground-4B-16C-DPO** | 62.59 | 2712 |
| UI-TARS-2B | 58.89 | 2687 | **Phi-Ground-4B-29C** | 61.04 | 2901 |
| UI-TARS-7B | 61.58 | 2745 | **Phi-Ground-7B-7C** | 59.25 | 2668 |
| UI-TARS-1.5-7B | 61.40 | 2950 | **Phi-Ground-7B-16C** | *63.02* | 2813 |
| UGround-7B | 52.96 | 4371 (hf) | **Phi-Ground-7B-16C-DPO** | **64.39** | 2813 |
| UGround-v1-7B | 57.99 | 2717 | **Phi-Ground-7B-29C** | 61.93 | 3103 |
| *Agent setting (Use long REs) with O4-mini as planner* | | | | | |
| seeclick | 19.60 | 8847 (hf) | **Phi-Ground-4B-7C** | 69.60 | 8122 |
| OS-ATLAS-4B | 17.99 | 9288 (hf) | **Phi-Ground-4B-16C** | 72.12 | 8212 |
| OS-ATLAS-7B | 45.50 | 9788 (hf) | **Phi-Ground-4B-16C-DPO** | *73.51* | 8212 |
| UI-TARS-2B | 66.37 | 8188 | **Phi-Ground-4B-29C** | 69.96 | 8401 |
| UI-TARS-7B | 69.78 | 8243 | **Phi-Ground-7B-7C** | 71.94 | 8168 |
| UI-TARS-1.5-7B | 71.58 | 8454 | **Phi-Ground-7B-16C** | 72.12 | 8313 |
| UGround-7B | 62.41 | 9871 (hf) | **Phi-Ground-7B-16C-DPO** | **73.87** | 8313 |
| UGround-v1-7B | 66.73 | 8214 | **Phi-Ground-7B-29C** | 71.40 | 8603 |

Table 16: Gold dataset evaluation results.

| model | Gold dataset | | | | | | | Gold-S |
|---|---|---|---|---|---|---|---|---|
| | PhotoShop | ClipChamp | Excel | PowerPoint | Word | Windows Setting | AVG. | |
| *End-to-end model setting (Use short REs)* | | | | | | | | |
| SeeClick-10B [44] | 1.41 | 14.53 | 21.50 | 13.41 | 11.63 | 59.76 | 20.37 | 51.66 |
| UGround-7B [30] | 29.11 | 65.36 | 48.60 | 58.54 | 41.86 | 86.01 | 54.91 | 74.88 |
| UGround-v1-7B [30] | 49.30 | **79.33** | 57.01 | 69.51 | 47.67 | 95.73 | 66.42 | 84.36 |
| OS-ATLAS-4B [31] | 5.63 | 21.23 | 14.95 | 13.41 | 16.28 | 60.55 | 22.01 | 47.87 |
| OS-ATLAS-7B [31] | 32.86 | 70.39 | 27.10 | 31.71 | 39.53 | 90.91 | 48.75 | 66.35 |
| UI-TARS-2B [28] | 46.95 | 67.60 | 55.14 | 57.32 | 48.84 | 84.90 | 60.12 | 79.15 |
| UI-TARS-7B [28] | 56.34 | *76.54* | 72.90 | 85.37 | 75.58 | 93.75 | 76.75 | **87.20** |
| UI-TARS-1.5-7B [29] | 63.85 | 65.92 | **83.18** | 80.49 | 81.40 | 88.62 | 77.24 | 86.73 |
| **Phi-Ground-4B-7C** | 46.01 | 62.57 | 52.34 | 65.85 | 58.14 | 91.86 | 62.8 | 76.78 |
| **Phi-Ground-4B-16C** | 58.22 | 73.18 | 76.64 | 81.71 | 76.74 | 93.28 | 76.63 | 85.78 |
| **Phi-Ground-4B-16C-DPO** | 64.79 | 68.72 | 78.50 | 87.80 | 75.58 | 93.99 | 78.23 | **87.20** |
| **Phi-Ground-4B-29C** | 56.81 | 68.72 | *80.37* | 87.80 | 73.26 | 94.15 | 76.85 | 83.41 |
| **Phi-Ground-7B-7C** | 57.28 | 69.83 | 66.36 | 81.71 | 68.60 | 91.94 | 72.62 | 85.31 |
| **Phi-Ground-7B-16C** | *69.95* | 69.83 | 77.57 | 87.8 | **77.91** | 93.28 | *79.39* | 84.83 |
| **Phi-Ground-7B-16C-DPO** | **72.30** | 69.27 | 72.90 | **90.24** | 77.91 | *95.18* | **79.63** | 84.36 |
| **Phi-Ground-7B-29C** | 67.61 | 67.60 | 79.44 | **90.24** | 72.09 | **96.28** | 78.88 | 84.83 |
| *Agent setting (Use long REs) with GPT-4O as planner* | | | | | | | | |
| SeeClick-10B | 1.41 | 9.50 | 11.21 | 6.10 | 2.33 | 48.46 | 13.17 | 32.23 |
| UGround-7B | 38.97 | 64.25 | 57.94 | 56.10 | 53.49 | 86.88 | 59.60 | 81.52 |
| UGround-v1-7B | 52.11 | 73.18 | 67.29 | 74.39 | 62.79 | 95.57 | 70.89 | 84.83 |
| OS-ATLAS-4B | 7.98 | 18.44 | 15.89 | 15.85 | 10.47 | 61.34 | 21.66 | 38.39 |
| OS-ATLAS-7B | 31.92 | 62.57 | 26.17 | 41.46 | 50.00 | 87.67 | 49.97 | 67.30 |
| UI-TARS-2B | 55.40 | 76.54 | 71.96 | 80.49 | 69.77 | 93.91 | 74.68 | 87.20 |
| UI-TARS-7B | 62.91 | 74.86 | 78.50 | 82.93 | 81.40 | 94.39 | 79.17 | 85.31 |
| UI-TARS-1.5-7B | 66.67 | 62.57 | 85.05 | 84.15 | **87.21** | 86.88 | 78.76 | 88.63 |
| **Phi-Ground-4B-7C** | 65.26 | 71.51 | 68.22 | 78.05 | 75.58 | 96.76 | 75.90 | 88.63 |
| **Phi-Ground-4B-16C** | 70.89 | **78.77** | 84.11 | 86.59 | 83.72 | 96.92 | 83.50 | **91.47** |
| **Phi-Ground-4B-16C-DPO** | **76.06** | *77.65* | *87.85* | 89.02 | 84.88 | *97.71* | **85.52** | *91.00* |
| **Phi-Ground-4B-29C** | 70.89 | 70.95 | 78.50 | **91.46** | 80.23 | 96.36 | 81.40 | 89.57 |
| **Phi-Ground-7B-7C** | 69.48 | 75.42 | 75.70 | 87.80 | 74.42 | 97.15 | 80.00 | 88.63 |
| **Phi-Ground-7B-16C** | 73.24 | 75.42 | **89.72** | **91.46** | 81.40 | 96.36 | 84.60 | 90.05 |
| **Phi-Ground-7B-16C-DPO** | **76.06** | 75.98 | 86.92 | 89.02 | *86.05* | **97.94** | *85.33* | *91.00* |
| **Phi-Ground-7B-29C** | 71.83 | 73.18 | 84.11 | 89.02 | 77.91 | 97.00 | 82.17 | 88.15 |
| *Agent setting (Use long REs) with O4-mini as planner* | | | | | | | | |
| SeeClick-10B | 2.35 | 20.11 | 9.43 | 4.94 | 6.90 | 49.72 | 15.57 | 39.42 |
| UGround-7B | 39.44 | 70.39 | 56.60 | 54.32 | 48.28 | 87.75 | 59.46 | 83.17 |
| UGround-v1-7B | 52.11 | 77.09 | 67.92 | 80.25 | 68.60 | 96.76 | 73.79 | 87.98 |
| OS-ATLAS-4B | 6.57 | 17.88 | 6.60 | 16.05 | 9.20 | 53.04 | 18.22 | 37.02 |
| OS-ATLAS-7B | 35.68 | 65.92 | 33.96 | 39.51 | 49.43 | 88.62 | 52.19 | 68.27 |
| UI-TARS-2B | 57.28 | 82.12 | 73.58 | 81.48 | 73.26 | 94.86 | 77.10 | 85.58 |
| UI-TARS-7B | 67.14 | 82.12 | 77.36 | 88.89 | 77.91 | 95.26 | 81.45 | 90.87 |
| UI-TARS-1.5-7B | 63.38 | 74.86 | 85.85 | 81.48 | **87.21** | 89.09 | 80.31 | 90.38 |
| **Phi-Ground-4B-7C** | 70.42 | 82.12 | 73.58 | 81.48 | 73.26 | 99.13 | 80.00 | 92.79 |
| **Phi-Ground-4B-16C** | 75.59 | 88.27 | 83.96 | 87.65 | 82.56 | 98.97 | 86.17 | **95.19** |
| **Phi-Ground-4B-16C-DPO** | *79.81* | 88.83 | **88.68** | 88.89 | *84.88* | *99.21* | **88.38** | **95.19** |
| **Phi-Ground-4B-29C** | 73.71 | 87.71 | 83.96 | 90.12 | 82.56 | 98.97 | 86.17 | 90.87 |
| **Phi-Ground-7B-7C** | 71.83 | 85.47 | 78.30 | 88.89 | 80.23 | **99.29** | 84.00 | 92.79 |
| **Phi-Ground-7B-16C** | 76.53 | *89.39* | *86.79* | 91.36 | 79.07 | 98.81 | 86.99 | 92.31 |
| **Phi-Ground-7B-16C-DPO** | **81.22** | **89.94** | 85.85 | 91.36 | 81.40 | *99.21* | *88.16* | 93.75 |
| **Phi-Ground-7B-29C** | 76.53 | 86.03 | 84.91 | 90.12 | 80.23 | 98.10 | 85.99 | 92.31 |

# C  Coordinates Representations and Loss

As discussed in the main text, we experimented with various coordinate representations and loss function designs. We found that these techniques can accelerate training when dealing with small datasets. However, when the training dataset exceeds 1 million samples, these methods do not exhibit significant improvements. Consequently, the content presented in this section highlights approaches that failed to scale. We disclose these findings to help future researchers avoid similar pitfalls.

Overall, the development of all the techniques discussed in this section stems from the following considerations: Unlike regression loss, modeling with natural language treats the difference between "19" and "20" as a gap of two tokens, which should theoretically be equivalent in distance to that between "18" and "19". Furthermore, differences in the units, tens, and hundreds places should have varying impacts—we might tolerate errors in the units place, but errors in the hundreds place are entirely unacceptable. These aspects highlight a disparity to regression loss and our expectations. However, experimental results indicate that using the most straightforward next token prediction and expressing coordinates in natural language is well enough, and there is no significant difference in outcomes among these techniques when the batch size is extremely large and the training volume is very high.

## C.1  Tokenized Coordinates

It has been observed that in most LLM, numbers are tokenized by digits, meaning the number 123 would be tokenized into three separate tokens: "1", "2", and "3". This form of representation offers limited interpretability in the context of images. In previous work, many researchers have modeled regions within an image using newly introduced special tokens. We also attempted this approach; however, when processing screenshots, the images often have extremely high resolutions, and buttons are relatively small. If we divide both the height and width into 1,000 discrete intervals and assign a new token to each square region, this would add 1 million tokens to the model, which is entirely impractical. Therefore, we opted to model the coordinate values from 1 to 1,000 as new tokens, for example $< p\ 123 >$ for value 123, using two tokens to represent a single position: $< p\ 123 >< p\ 456 >$. In this way, we added only 1000 special tokens to the model. We discussed various strategies for initializing these 1000 tokens. To better illustrate our point, we first introduce the following definition: let the function $emb : V \to R^n$ represent the retrieval of the embedding for a specific token $v \in V$ from a pre-trained model. Then we define the following variables:

$$C_{rand} = random(mean = MEAN(\{emb(v), v \in V\}), std = STD(\{emb(v), v \in V\}))$$

$$C_{digit} = random(mean = MEAN(\{emb(v), v \in digit\}), std = STD(\{emb(v), v \in digit\}))$$

$$M_{digit} = \frac{1}{|digit|} \sum_{v \in dight} emb(v) \qquad E_p = emb(\text{"point"})$$

We then consider the following five initialization methods:

- The implementation in the Hugging Face Transformers library (hf) calculates the mean and variance of all pre-trained embeddings to generate a random vector, which is then used to initialize all newly added tokens.
- We adapted the hf method by restricting the embeddings used for calculating the mean and variance to only digit tokens, naming this approach R-digit.
- The digit-mean method directly uses the mean of the embeddings of digit tokens from the pre-trained model to initialize larger number tokens.
- The term main-digit refers to using the embedding of the digit in the hundredths place as the initialization. For example, for the number 234, the initialization would use the embedding of the digit "2."
- Prefix-learning method first freeze all parameters except the embedding of newly added tokens, and then used the learned embedding to initialize in the later training.

Additionally, we handle <point>, </point> and digits 0 to 9 differently. The specific assignments can be found in Table 17, where the training results are also presented.

We set the training volume for these experiments to 5 million samples and detailed setting can be found in Table 11. Regrettably, the results in Table 17 indicate that all initialization methods significantly underperform compared to directly using natural language to express coordinates. Additionally, during training, we observed slow convergence and significant fluctuations in the gradient norm. These phenomena suggest that the introduction of too many special tokens, which have not undergone large-scale pre-training, can interact adversely with the pre-trained parameters, leading to training collapse. However, when high resolution is required, 1000 special tokens become necessary. Thus, in the field of UI grounding, this technique appears to be less practical.

Table 17: The details and result of different initialization methods for newly added special tokens.

| method name | Initial value for special tokens | | | | | | | | | | Gold-S |
| | <point> | </point> | <p 0> | <p 1> | ... | <p 9> | <p 10> | <p 11> | ... | <p 999> | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| hf | $C_{rand}$ | $C_{rand}$ | $C_{rand}$ | $C_{rand}$ | | $C_{rand}$ | $C_{rand}$ | $C_{rand}$ | | $C_{rand}$ | 81.7 |
| R-digit | $E_p$ | $E_p$ | $C_{digit}$ | $C_{digit}$ | | $C_{digit}$ | $C_{digit}$ | $C_{digit}$ | | $C_{digit}$ | 83.2 |
| digit-mean | $E_p$ | $E_p$ | $emb("0")$ | $emb("0")$ | | $emb("9")$ | $M_{digit}$ | $M_{digit}$ | | $M_{digit}$ | 74.6 |
| main-digit | $E_p$ | $E_p$ | $emb("0")$ | $emb("0")$ | | $emb("0")$ | $emb("0")$ | $emb("0")$ | | $emb("9")$ | 12.6 |
| prefix-learning | - | - | - | - | | - | - | - | | - | 82.3 |
| Natural language | - | - | - | - | | - | - | - | | - | 88.9 |

## C.2 Label Smoothing

It is fairly intuitive to consider that applying label smoothing to digit tokens might produce an effect similar to that of regression loss. Specifically, for instance, when the ground truth token is the digit 5, the standard Cross-Entropy loss assigns a label of 1 to digit 5 and 0 to all other tokens. However, if we apply label smoothing and assign a certain degree of smoothing to digits 4 and 6, we effectively inform the model that digits 5 and 4, 6 are numerically close. This approach results in an outcome akin to regression loss. In fact, it is possible to derive how to set label smoothing parameters to achieve equivalence with using regression loss.

We will next derive the formula for our label smoothing technique. First, we provide the following notations. Let $\mathbf{x} = (x_0, x_1, ...x_n)$ be a tokenized sentence and $x_n$ is a digit token whose digit value(integer number) is $T$, $t$ is the token id of $x_n$. We usually model the probability given by language model as $p(\mathbf{x}) = p(x_0)p(x_1|x_0)\dots p(x_n|x_0, x_1, \dots, x_{n-1})$. Define $V$ as the vocabulary set, $V_d$ as the set of digit token and $V = V_d \cup V_t$. Now let's recall the formulation of classical language modeling loss:

$$L_{lm} = -(\log p(x_0) + \log p(x_1|x_0) + \dots + \log p(x_n|x_0, \dots, x_{n-1}))$$

The last term:

$$L_n = -\log p(x_n|x_0, \dots, x_{n-1}) := -\log p_n = -\sum_{i=1}^{|V|} y^{(i)} \log p^{(i)},$$

where $p^{(i)}$ is the softmax result of last hidden activations and $y_i$ is the label. In general cross-entropy loss, the label is just a one-hot encoding:

$$y^{(i)} = \begin{cases} 1 & \text{if } i = t \\ 0 & \text{else} \end{cases}$$

In the following, we will derive how to approximate the regression loss with an appropriate designed label $y^{(i)}$. For $k \in V_d$ let the digit number of digit token id $k$ be $K$. We design a regularized MSE loss as:

$$L_{MSE} = \mathbb{E}_{k \in V_d}[(K - T)^2] + \psi \mathbb{E}_{k \in V_t}[1],$$

where the second term is a punishment for other none-digit tokens and $\psi$ is the punishment factor. We can calculate that:

$$
\begin{aligned}
L_{MSE} &= \mathbb{E}_{k \in V_d}[(K-T)^2] + \psi \mathbb{E}_{k \in V_t}[1] \\
&= \sum_{k \in V_d} (K-T)^2 p^{(k)} + \psi \sum_{k \in V_t} p^{(k)} \\
&= \sum_{k \in V_d/\{t\}} (K-T)^2 p^{(k)} - \psi \sum_{k \in V_d/\{t\}} p^{(k)} + \psi \sum_{k \in V_d/\{t\}} p^{(k)} + \psi \sum_{k \in V_t} p^{(k)} \\
&= \sum_{k \in V_d/\{t\}} [(K-T)^2 - \psi]p^{(k)} + \psi(1 - p^{(t)}) \\
&= -\psi \left\{ p^{(t)} + \sum_{k \in V_d/\{t\}} [1 - \frac{(K-T)^2}{\psi}]p^{(k)} - 1 \right\}
\end{aligned}
$$

This will be equivalent to optimize loss: $\tilde{L}_{MSE} = -\sum_{i=1}^{V} \tilde{y}^{(i)} p^{(i)}$, where:

$$\tilde{y}^{(i)} = \begin{cases} 1 & \text{if } i = t \\ 1 - \frac{d(K,T)}{\psi} & \text{if } i \in V_d \setminus \{t\} \\ 0 & i \in V_t \end{cases} \tag{1}$$

$d(K,T)$ is the distance function, for MSE, it should be $d(K,T) = (K-T)^2$. We directly use Equation 1 for label smoothing, denoted as $\hat{L}_{MSE} = -\sum_{i=1}^{V} \tilde{y}^{(i)} \log p^{(i)}$. It is important to note that this differs from $\tilde{L}_{MSE}$. However, by leveraging the convexity of the logarithmic function, we can easily demonstrate that $\hat{L}_{MSE}$ serves as an upper confidence bound for $\tilde{L}_{MSE}$. From a practical standpoint, the information embedded in Equation 1 suggests that the label values are larger for positions closer to the target, which is highly intuitive.

Table 18: The scaling effect of label smoothing technique is very limited.

| $d(K,T)$ | $\psi$ | $N = 50K, \quad BS = 64$ ScreenSpot-V2 | Gold-S | $N = 1M \quad BS = 2048$ ScreenSpot-V2 | Gold-S |
|---|---|---|---|---|---|
| $d(K,T) = (K-T)^2$ | 10 | 52.5 | 66.9 | 81.7 | 85.9 |
| $d(K,T) = (K-T)^2$ | 30 | 55.8 | 66.2 | 83.3 | 88.1 |
| $d(K,T) = |K-T|$ | 10 | 54.6 | 67.6 | **84.3** | 87.8 |
| $d(K,T) = |K-T|$ | 30 | **60.1** | **69.6** | 84.1 | 87.6 |
| no label smoothing | | 52.3 | 63.2 | 84.2 | **88.8** |

The results, as shown in Table 18, indicate that this technique demonstrates a clear advantage when dealing with smaller datasets and smaller batch sizes. However, when we increased the training data size to 1 million and the batch size to 2048, we observed little to no improvement, with results falling within the margin of fluctuation. This suggests that while the approach may offer some acceleration benefits in resource-constrained scenarios, it holds limited significance for large-scale training. In such cases, with larger batch sizes, the model's optimization direction aligns with the regression loss, reducing the scaling advantage of this technique.

## C.3 Loss Re-weighting

We attempt to assign different weights to the loss of different tokens. For instance, we assign higher weights to tokens in the tens place compared to those in the units place. This approach ensures that the model prioritizes the correctness of the tens place over the units place. If expressed formally using equations, it can be represented as follows:

$$L_{re} = -\sum_{t=0}^{n} w_t \log p(x_t | x_{t-1}..., x_0).$$

Initially, we set $w_t = 1.0$ when $x_t$ is not a digit token. Then we consider different settings of weights for digit tokens, as shown in Table 19.

Table 19: Experiment on selecting parameters for loss reweighting.

| weights for digit position hundrads | tens | units | $N = 50K, \quad BS = 64$ ScreenSpot-V2 | Gold-S | $N = 1M \quad BS = 2048$ ScreenSpot-V2 | Gold-S |
|---|---|---|---|---|---|---|
| 1.0 | 1.0 | 1.0 | 52.3 | 63.2 | 84.2 | **88.8** |
| 2.0 | 1.5 | 1.0 | 6.3 | 16.2 | 15.5 | 22.7 |
| 4.0 | 2.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 1.0 | $1/10$ | $1/100$ | 55.4 | **66.8** | **84.5** | 85.8 |
| 1.0 | $1/\sqrt{10}$ | $1/10$ | 54.7 | 63.6 | 83.3 | 86.7 |
| 1.0 | $1/\ln 10$ | $1/(\ln 10)^2$ | **57.6** | 64.8 | 84.1 | 86.7 |

The first block in the table represents the control group without using reweighting techniques. We first confirm that the weight of digit tokens cannot exceed 1.0, even though these format-related tokens appear in almost all data. We found that if the weight of digit tokens is even slightly greater than 1.0, it causes the model to output in an unexpected format

during testing. This results in parsing errors, causing the test results to be nearly zero. In contrast, in models trained under normal conditions, the proportion of parsing errors is nearly zero.

When the weight of digit tokens is proportionally adjusted to smaller than 1.0 and the model is trained accordingly, we obtain results similar to those described in Section C.2. Specifically, when training resources are extremely limited, we observe that the reweighting technique accelerates model convergence and achieves consistently better results in low-sample scenarios. However, when the data size and batch size increase significantly, the benefits of this technique are minimal or unstable. Such results are insufficient to support us to widespread this technique. The reason for this phenomenon may be that when the training volume and batch size increase, learning based on higher numeric values exhibits greater certainty and stability (or lower perplexity), which facilitates more effective learning. In contrast, learning based on lower numeric values might fluctuate due to errors present in the dataset, leading to slower learning. This process is similar to our reweighting technique, and therefore, as the training volume and batch size increase, this technique is effectively replaced.

## D    Data Pre-processing Details

### D.1    CommonCrawl Data Pre-processing

#### D.1.1    Rendering Resolutions and Filtering

**Rendering resolutions.**    When rendering each web page, we initially select a screen size with equal probability (1/3 chance) from the options of 1080p, 2.5k, and 4k. The corresponding screen area ($Space$) are $1920 \times 1080$, $2560 \times 1440$, and $3840 \times 2160$, respectively. For a given screen area, we then randomly choose an aspect ratio ($R_w, R_h$) from the following set of aspect ratios:

$$(R_w, R_h) \in \left\{ (1 + \frac{i}{N}, 2 - \frac{i}{N}) | i = 0, 1, ..., N \right\}.$$

The final screen size $(W, H)$ used for rendering can be calculated using the following formula:

$$(W, H) = (R_w \times S, R_w \times S), \qquad S = \left\lceil \sqrt{\frac{Space}{R_w R_h}} \right\rceil$$

**Filtering.**    During webpage rendering, we implemented certain filtering processes. Unlike the filtering described in Section D.1.2, this stage of filtering requires the relevant code (such as JavaScript) used in browser rendering. In contrast, the subsequent filtering refers to offline filtering conducted after the necessary information has been stored. Therefore, even though both processes involve rule-based filtering, they are described in two separate sections. During webpage rendering, we can use JavaScript to obtain interactive information, and based on this information, we have established the following filtering rules for all HTML elements.

We first retain only those elements that meet any of the following conditions:

- Interactive Tags: The HTML tag name of the element is one of *'button'*, *'input'*, *'textarea'*, *'select'*, *'a'*, *'form'*
- Event Attribute: The elements have specific JavaScript methods (functions) attached, such as *'onclick'*, *'onmousedown'*, *'onmouseup'*, *'onmouseover'*, *'onmouseout'*, *'onkeydown'*
- Role attribute: Elements with the following role attributes are generally interactive: *'button'*, *'link'*, *'textbox'*, *'menuitem'*, *'option'*, *'checkbox'*, *'radio'*, *'tab'*, *'switch'*.
- Interactive class: The class name of the element is a string type and the class name is one of *'btn'*, *'button'*, *'input'*, *'link'*, *'nav'*, *'menu'*, *'item'*.
- Is icon: Tag name is one of *'i'*, *'span'*, *'svg'* and the class name is one of *'fa'*, *'fas'*, *'far'*, *'fal'*, *'fab'*, *'material-icons'*
- Is image: The tag name is *'img'*

Subsequently, we remove all elements that are not visible on the screen, such as those that require scrolling to be viewed. We then store all relevant information of the elements that meet the criteria, in JSON format, along with the rendered images. This includes details such as bounding box coordinates, HTML code, and various attributes. We also store a layout diagram, which matches the size of the rendered screenshot. However, different types of elements are represented by different colors occupying their respective areas; for example, interactive text is marked in red, and images in cyan. This type of layout diagram allows for the expression of button positions without focusing on the content itself and is used for data deduplication at the webpage level.

### D.1.2 Offline Rule-based Filtering

Once the data has been stored, we consider the following filtering rules:

**Boxes deduplication.** If a box completely encompasses multiple other boxes, we first remove the outer box. Such boxes are typically div containers of a module area, containing multiple related buttons. When one box contains another (determined by an IoU greater than a certain threshold), we remove the larger box. This situation is common in web design with nested containers and errors generated by OmniParser when using it to create boxes.

**Remove empty boxes.** In both webpage rendering and OmniParser, there are instances where certain boxes appear in a blank, solid-colored area devoid of any content. For each candidate box, we crop the corresponding region from the screenshot and use the pixel standard deviation to directly determine if the area is solid-colored. If it is deemed to be solid-colored, we delete the box.

**Text content recognition.** Sometimes, the text content within a screenshot, such as a long sentence, is recognized as an element. We wish to retain buttons with text but not these non-interactive content texts. To achieve this, we use the aspect ratio of the box as a filtering criterion. If the aspect ratio of a box exceeds a certain threshold, we discard that box.

### D.2 Re-sampling algorithm

---
**Algorithm 3** Re-sampling algorithm
---
**Require:** The center point set of training dataset $C = \{(x_i, y_i)\}$ using relative coordinates. Segmentation granularity $N, M$, Sampling factor $\psi$.

**Ensure:** Sampled center point set $\hat{C}$

1: all_box $\leftarrow \{(i,j) : \text{list}() \mid i \in \{0, 1, ...N-1\}, j \in \{0, 1, ...M-1\}\}$
2: **for** $(x, y)$ **in** $C$ **do**
3:      locate $\leftarrow (\text{int}(x//(1.0/N)), \text{int}(y//(1.0/M)))$
4:      all_box[locate].append$((x, y))$
5: **end for**
6: dist $\leftarrow [\text{len}(v) \mid v \in \text{all\_box.values}()]$
7: dist.sort()
8: keep_number $\leftarrow$ dist$[\text{int}(N \times M \times \psi)]$
9: $\hat{C} \leftarrow \text{list}()$
10: **for** v **in** all_box **do**
11:      N_Sample $\leftarrow$ **min**$(\text{len}(v), \text{keep\_number})$
12:      $\hat{C}$.extend(**random_sample**$(v, \text{N\_Sample})$)
13: **end for**
---

# E  Prompts

## E.1  Reference expression generation prompt.

---

**Prompt for generating Long-Gold RE / Training data generation**

System prompt:

---

User will provide you with a screenshot, in which a specific area will be highlighted with a red rectangular box. We will also provide a cropped image of the corresponding area, and (optionally) additional information related to the area to help you understand it. Your task is to generate several references regarding the target area on the original image.

Specific task requirements are as follows: You will analyze and output a dictionary in JSON file format. The key-value pairs included are as follows:

- area_type: Choose one from 'icon', 'text'. These represent whether the target area is an indicative icon, text.

- interactive: bool, indicating whether this element in the screenshot scenario is an interactive element (e.g., clickable, inputable, etc.). If it is static text or an image, then it is not interactive.

- context: Generate a background context describing why the current screen and area would be used. For example, if the area is the close button of a image file, the context could be that the user is editing a file and has completed their task at this moment.

- functional_reference: A reference about the target area, involving the function of the target area.

- positional_reference: A reference about the target area by describing the position of the target area, such as layout and nearby elements.

- appearance_reference: A reference about the target area by describing the appearance of the target area.

**Ensure that anyone can uniquely identify this area in the screenshot through any one of the references. Don't mention red rectangular box.**
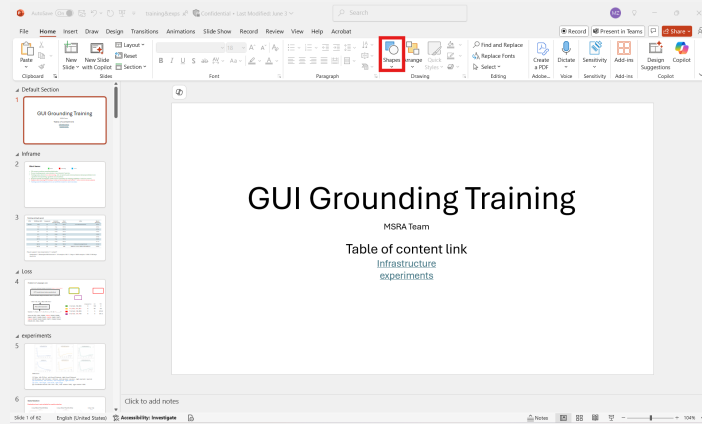**Your output references should only include the element description itself and follow the requirements. Do not start with "the target element" or "the element"**
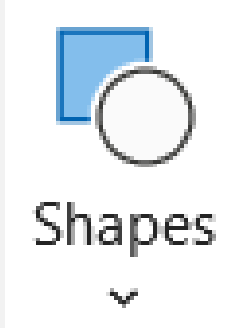Your output should follow this format strictly:

```
# Analyze
A free form analyze of the screen shot and task.
# Output
```json
{
    "area_type": ...,
    "interactive": ...,
    "context": ...,
    "functional_reference": ...,
    ...,
    "appearance_reference": ...
}
```
```

---

Model input:

---

# Screenshot with highlight



# Cropped target image



Model output:

---

# Analyze

The screenshot shows a PowerPoint presentation interface. The highlighted area is located on the toolbar of the PowerPoint application and is labeled as "Shapes." This area is used to insert various shapes into the presentation slides, an essential feature for users looking to enhance their slides with graphical elements. The toolbar section where the "Shapes" button is located is a common area for formatting and editing tools, making it easily accessible for users during the slide creation process.

# Output

```
{
    "area_type": "icon",
    "interactive": true,
    "context": "While working on a PowerPoint presentation, the user is
    likely looking to add graphical elements or illustrations to their
    slides to enhance the visual appeal or to convey information more
    effectively.",
    "functional_reference": "Shapes button is used to add various
    graphical shapes to a slide in a PowerPoint presentation.",
    "positional_reference": "Located on the toolbar under the 'Home' tab,
    to the right of the 'Layout' button and to the left of the 'Arrange'
    button.",
    "appearance_reference": "Contains a white circle and a blue square,
    with the label 'Shapes' underneath."
}
```

### E.2  Planner's generation prompt for agent setting.

---

**Prompt for generating Long RE / Agent settings for evaluation**

System prompt:

---

User will provide you with a screenshot and a short instruction related to the area to help you understand it.Your task is to classify the target area and generate several references regarding the target area on the original image. Specific task requirements are as follows: You will analyze and output a dictionary in JSON file format. The key-value pairs included are as follows:

- context: Generate a background context describing why the current screen and area would be used. For example, if the area is the close button of a image file, the context could be that the user is editing a file and has completed their task at this moment.
- functional_reference: A reference about the target area, involving the function of the target area.
- positional_reference: A reference about the target area by describing the position of the target area, such as layout and nearby elements.
- appearance_reference: A reference about the target area by describing the appearance of the target area.
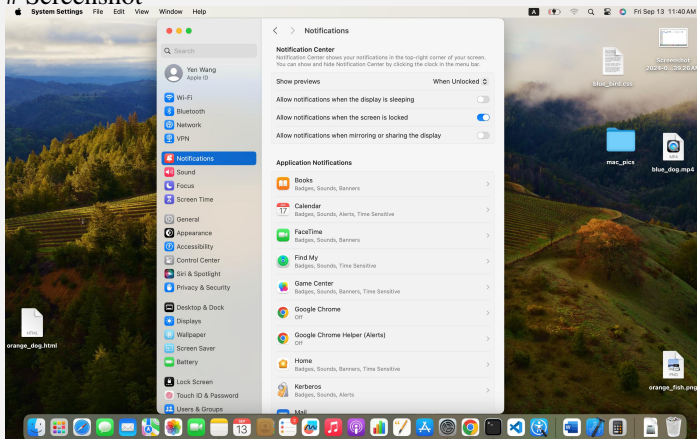
**Ensure that anyone can uniquely identify this area in the screenshot through anyone of the references. Ensure that the reference is complete and independent.**
Your output should follow this format strictly:

```
# Analyze
A free form analyze of the screen shot and task.
# Output
```json
{
    "context": ...,
    "functional_reference": ...,
    ...,
    "appearance_reference": ...
}
```
```

Model input:

---

# Screenshot



# Instruction
adjust the options of "show previous"

---

Model output:

---

# Analyze
The screenshot displays a macOS System Settings window focused on the "Notifications" section. The user intends to adjust the "Show previews" option, which is part of the Notification Center settings. This area allows users to control how notifications are displayed on their device, particularly when the screen is unlocked. The screen is divided into two main sections: the left pane with different settings categories and the right pane showing detailed options for the selected category, which is "Notifications" in this case.
# Output

```
{
    "context": "The user is configuring notification settings to determine
    how previews of notifications are displayed on their device when the
    screen is unlocked.",
    "functional_reference": "The 'Show previews' option allows the user to
    control whether notification previews are shown on the screen when it
    is unlocked.",
    "positional_reference": "The 'Show previews' option is located in the
    right pane of the System Settings window under 'Notification Center'
    settings, just below the descriptive text about notifications.",
    "appearance_reference": "The 'Show previews' option is a dropdown menu
    located next to the text 'Show previews' with the current setting
    'When Unlocked'."
}
```
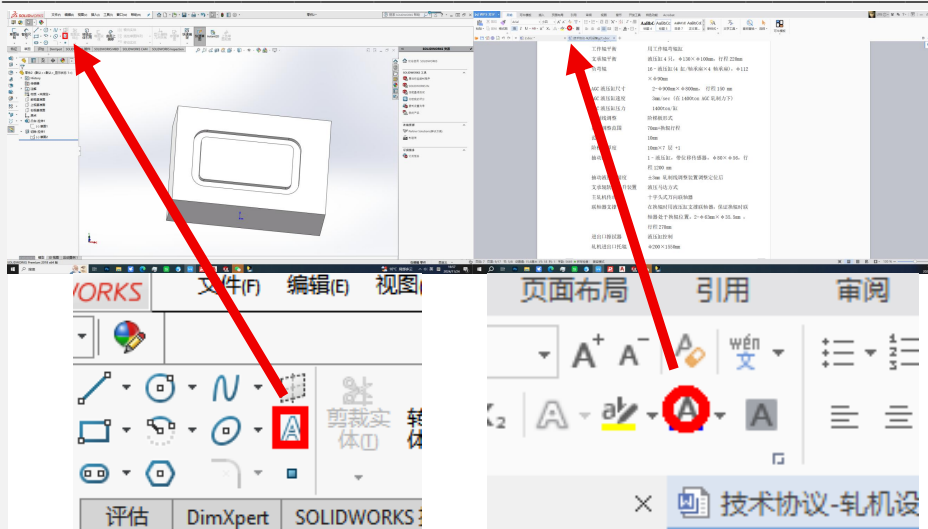
# F   More Cases Study

---

Error grounding case study with human level reference. Case-1: Similar Icons

Reference:

---

Toggles visibility of all annotations (dimensions, notes, symbols) in the graphics area. Located in the heads-up view toolbar at the top center of the viewport, immediately to the right of the Temporary Axes button and left of the Measure button. Blue uppercase letter A on a light gray square background.

---

Output:

---
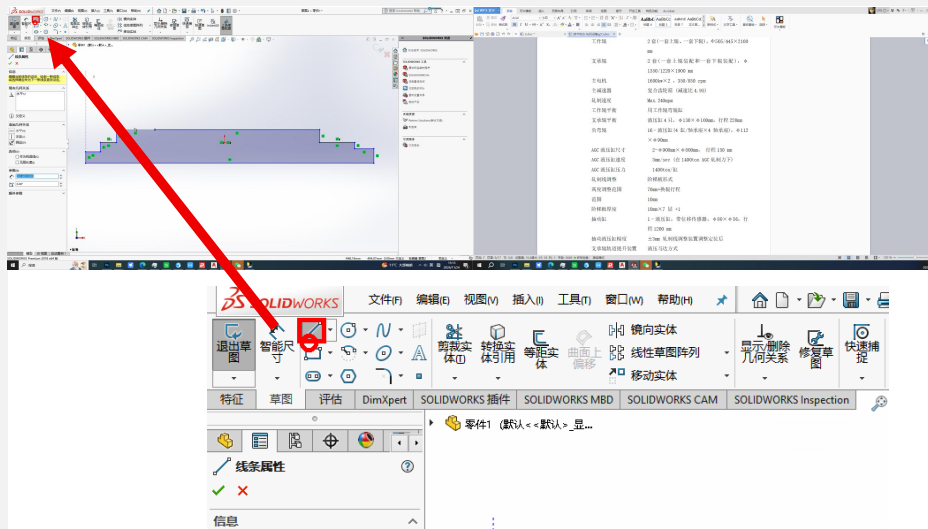
Error grounding case study with human level reference. Case-2: Precision Error

Reference:

Convert Entities command projects selected edges or curves from the model into the active sketch as sketch entities.Located on the Quick Access Toolbar at the very top of the SOLIDWORKS window, immediately to the right of the Redo icon and before the application menu bar.Grey square button showing a diagonal blue dashed line connecting two white square endpoints.
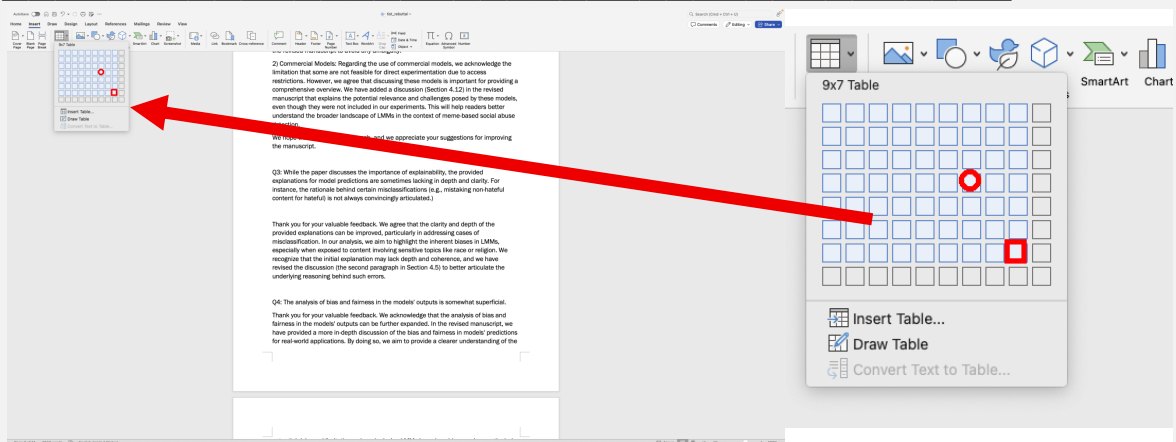
Output:



Error grounding case study with human level reference. Case-3: Lack spatial reasoning

Reference:

Grid cell for selecting a 9x7 table dimension when inserting a new table. Cell in the ninth column of the seventh row within the table size preview grid under the Table button in the Insert tab. Light-blue interior square outlined by a white inner border and a thicker blue outer border.
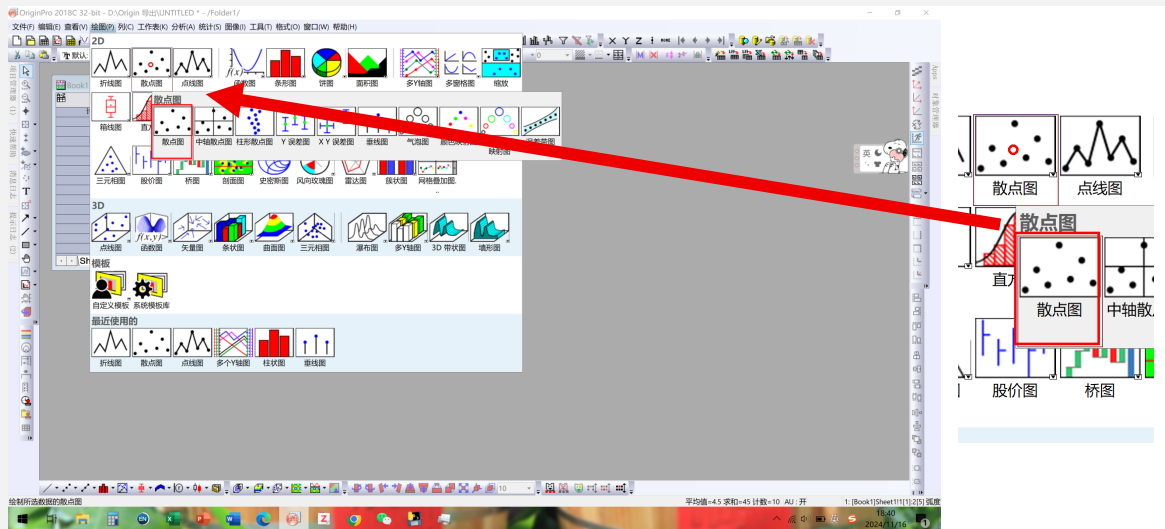
Output:

Error grounding case study with human level reference. Case-4: Two same area

Reference:

scatter plot template used to generate an X-Y scatter diagram from worksheet data. icon located in the top section of the 2D plot palette, in the row of scatter type graphs beneath the line and bar icons; it is the second icon from the left in that row. white rectangular button with several solid black dots arranged randomly inside and the Chinese label displayed underneath

Output:



Error grounding case study with human level reference. Case-5: Interactive area

Reference:

Toggles the applied Digital Glitch effect on or off for live comparison in the Effect Controls panel. Found immediately to the left of the "Digital Glitch" effect name under the adjustment layer's effects list in the top-left panel. Small grey lowercase "fx" icon on a dark background matching the style of other effect-toggle buttons.

Output: