



SEEClick: HARNESSING GUI GROUNDING FOR ADVANCED VISUAL GUI AGENTS

Kanzhi Cheng^{◇◇*} Qiushi Sun[♡] Yougang Chu[◇] Fangzhi Xu[♡]
 Yantao Li[◇] Jianbing Zhang[◇] Zhiyong Wu[♡]

[◇]Department of Computer Science and Technology, Nanjing University

[♡]Shanghai AI Laboratory

{chengkz, chuyg, li_yantao}@smail.nju.edu.cn qiushisun@u.nus.edu
 fangzhixu98@gmail.com zjb@nju.edu.cn wuzhiyong@pjlab.org.cn

ABSTRACT

Graphical User Interface (GUI) agents are designed to automate complex tasks on digital devices, such as smartphones and desktops. Most existing GUI agents interact with the environment through extracted structured data, which can be notably lengthy (e.g., HTML) and occasionally inaccessible (e.g., on desktops). To alleviate this issue, we propose a visual GUI agent – SeeClick, which only relies on screenshots for task automation. In our preliminary study, we have discovered a key challenge in developing visual GUI agents: GUI grounding – the capacity to accurately locate screen elements based on instructions. To tackle this challenge, we propose to enhance SeeClick with GUI grounding pre-training and devise a method to automate the curation of GUI grounding data. Along with the efforts above, we have also created *ScreenSpot*, the first realistic GUI grounding benchmark that encompasses mobile, desktop, and web environments. After pre-training, SeeClick demonstrates significant improvement in *ScreenSpot* over various baselines. Moreover, comprehensive evaluations on three widely used benchmarks consistently support our finding that advancements in GUI grounding directly correlate with enhanced performance in downstream GUI agent tasks. The model, data and code are available at <https://github.com/njucckevin/SeeClick>.

1 INTRODUCTION

A perennial topic in machine intelligence is the development of Graphical User Interface (GUI) agent systems, like Siri and Copilot, to automate complex tasks on computing devices, thereby reducing human workload (Shi et al., 2017; Li et al., 2020a). Recent advances in Large Language Models (LLMs) such as GPT-4 (OpenAI, 2023) and LLaMA (Touvron et al., 2023) have significantly propelled the evolution of GUI agents (Gur et al., 2023; Zhou et al., 2023). These agents interact with the environment by interpreting structured texts, e.g., HTML from webpages, then elicit LLM for planning, reasoning, and execution (Kim et al., 2023; Zheng et al., 2023).

However, GUI agents depend on structured text face three inherent limitations: (1) Structured text is not always accessible, especially for iOS or desktop applications where acquiring such information is challenging (Shaw et al., 2023); (2) The verbose nature of structured text serves as an inefficient context for LLMs, while also omitting crucial information such as layout, images, and icons (Deng et al., 2023); (3) The variety of structured text - including HTML, DOM, and Android VH - necessitates the curation of task-specific observation and action spaces (Kim et al., 2023; Zhou et al., 2023). These entrenched deficiencies in text-based approaches call for an alternative solution.

In this paper, we introduce SeeClick, a visual GUI agent built on Large Vision-Language Models (LVLMs). Inspired by human interaction with GUIs, as illustrated in Figure 1, SeeClick is designed to perform low-level actions like clicking or typing directly by observing interface screenshots. This in-

* Work done during internship at Shanghai AI Laboratory.

novative approach bypasses the interaction with cumbersome structured text, empowering SeeClick to universally adapt to various GUI platforms. Building such visual agents involves a foundational challenge: GUI grounding - the capacity to accurately locate screen elements based on instructions, which is absent in current LVLMs. To tackle this challenge, SeeClick enhances LVLm with a GUI grounding pre-training strategy. We devise a method to automate the curation of web grounding data and adapt public mobile UI datasets to obtain mobile grounding data. SeeClick employs the above-curated dataset for continual pre-training of the LVLm, enabling it to accurately locate elements such as text, widgets, and icons in various GUI environments.

Given GUI grounding is a fundamental yet underexplored capacity for GUI agents, we created *ScreenSpot*, the first realistic GUI grounding evaluation benchmark across various GUI platforms. *ScreenSpot* contains over 600 screenshots and 1200 instructions from iOS, Android, macOS, Windows, and webpages, and specifically includes both text-based elements and a variety of widgets and icons. Evaluation results confirm SeeClick’s superiority over current LVLms, validating the effectiveness of GUI grounding pre-training.

Finally, we adapt SeeClick to mobile and web agent tasks, including MiniWob (Shi et al., 2017), AITW (Rawles et al., 2023), and Mind2Web (Deng et al., 2023). As a purely vision-based agent, SeeClick achieves impressive performance. It is noteworthy that SeeClick surpasses the strong visual baseline Pix2Act, utilizing merely 0.3% training data. Moreover, experimental results on three tasks consistently support our finding that improvement in GUI grounding directly correlates with enhanced agent task performance.

This paper makes the following contributions:

- We develop a unified visual GUI agent SeeClick, which performs clicking and typing actions rely solely on interface screenshots across diverse GUI platforms.
- We prospectively explore GUI grounding for visual GUI agents, and enhanced SeeClick with proposed GUI grounding pre-training strategy.
- We create a realistic GUI grounding benchmark *ScreenSpot*, encompassing more than 1200 instructions from various GUI platforms.
- Experimental results on *ScreenSpot* and three agent tasks demonstrate that enhancing agents’ grounding capacity is key to improving performance in downstream agent tasks.

2 RELATED WORK

Autonomous GUI Navigation Early research explored GUI task automation in simplified web (Shi et al., 2017; Liu et al., 2018; Gur et al., 2018) and mobile UI (Li et al., 2020a; Burns et al., 2022; Li & Li, 2022). With the recent advances in LLMs (OpenAI, 2023; Touvron et al., 2023; Xu et al., 2023; Sun et al., 2023; Kim et al., 2023), LLM-centric agents have become the dominant paradigm. A line of works focused on prompting ChatGPT and GPT-4 for web tasks, via in-context learning (Zheng et al., 2023) and self-refine (Kim et al., 2023). Other research explored training LLMs as specialized GUI agents. Deng et al. (2023) devised a two-stage method for identifying target elements within intricate HTML. Gur et al. (2023) proposed to interact with websites via synthesizing programs.

Given the constraints of LLM to only process text, recent efforts have attempted vision-based GUI navigation (Shaw et al., 2023; Zhan & Zhang, 2023; Hong et al., 2023). These methods primarily utilize GPT-4V (Yan et al., 2023; Gao et al., 2023) and also require GUI metadata as input (Yang et al., 2023a; Zheng et al., 2024). In this paper, we construct a universal visual GUI agent SeeClick

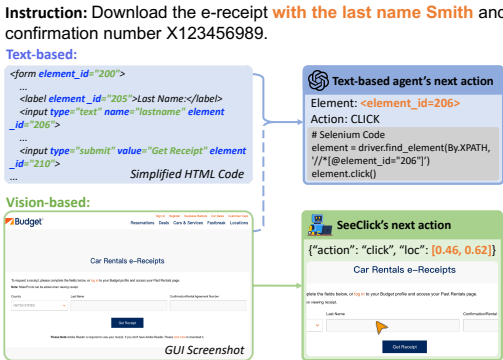
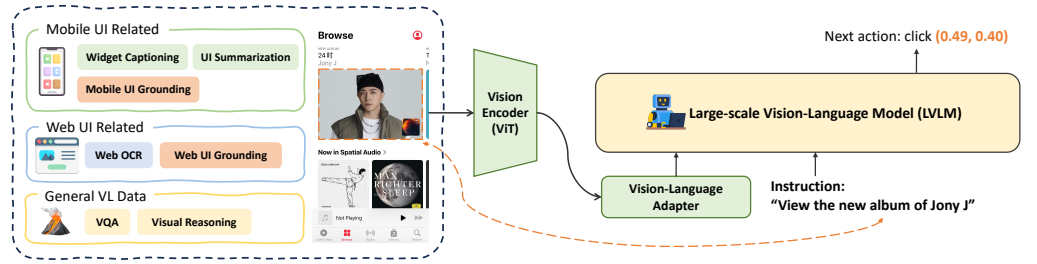
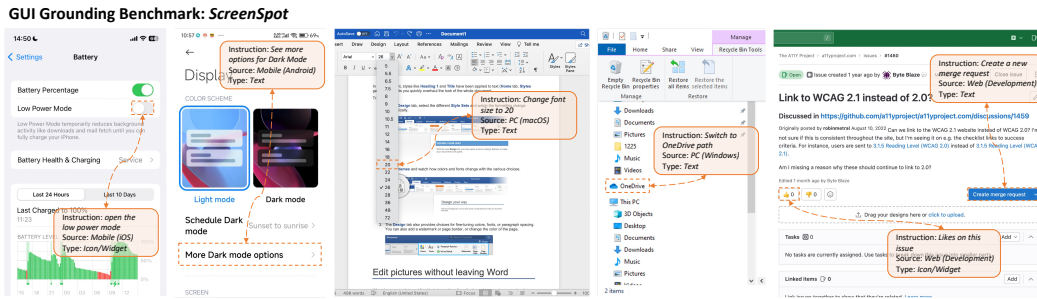


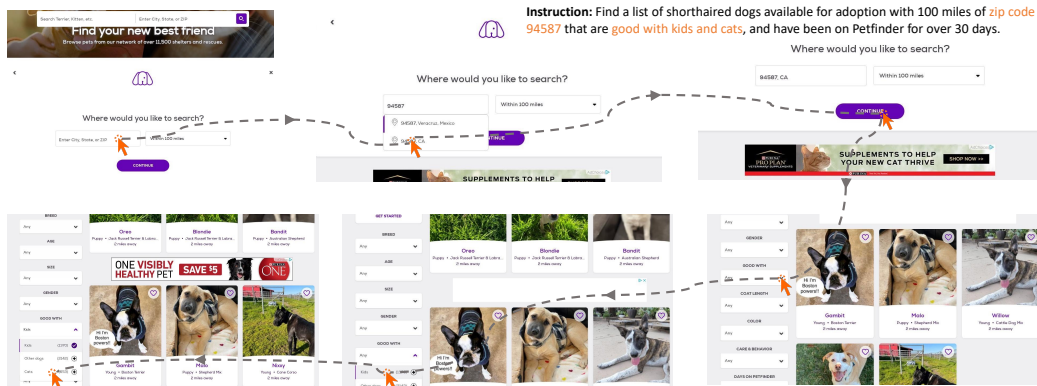
Figure 1: Text-based agents select target elements from structured texts, occasionally augmented with screenshots. SeeClick employs a vision-based methodology to predict action locations solely relying on screenshots.



(a) Overview of SeeClick's framework and GUI grounding pre-training.



(b) Examples of the proposed GUI grounding benchmark *ScreenSpot*.



(c) SeeClick as a visual GUI agent in downstream task.

Figure 2: Overview of our universal visual GUI agent SeeClick. (a) depicts the framework of SeeClick and GUI grounding pre-training. (b) provides examples of *ScreenSpot* across various GUIs and types of instructions. (c) displays the real-world application of SeeClick on downstream web agent task.

by customizing open-sourced LVLML, capable of operating across various GUI platforms without needing any GUI metadata.

Large Vision-Language Models Recently, researchers have invested tremendous effort in constructing LVLMLs capable of jointly processing image and text (Liu et al., 2023a; Zhu et al., 2023; Ye et al., 2023; Li et al., 2023). LVLMLs integrate vision encoders like ViT (Dosovitskiy et al., 2020), with LLMs such as LLaMA (Touvron et al., 2023), inheriting LLMs’ language and reasoning abilities to perform diverse vision-language tasks.

A series of studies focused on the grounding capabilities of LVLMLs (Wang et al., 2023; Bai et al., 2023; Chen et al., 2023a), such as providing bounding boxes for objects while generating responses (Chen et al., 2023b; Peng et al., 2023). Nonetheless, these efforts primarily addressed natural images and did not explore GUI contexts. This paper focuses on grounding in GUIs and explores the potential of LVLMLs as GUI agents.

3 APPROACH

In our preliminary study, we have discovered a key challenge in developing visual GUI agents: GUI grounding, the capacity to locate screen elements based on instructions. Although recent LVLMs have claimed visual grounding capability on natural images, GUI screenshots differ significantly with dense text and numerous icons and widgets. These differences impair existing LVLMs’ grounding performance in GUI contexts and limit their potential as visual GUI agents.

This paper seeks to harness LVLMs with GUI grounding skills, paving the way for a visual GUI agent that executes instructions only relying on screenshots. As presented in Figure 2, SeeClick is a foundational model for GUIs, and we adapt it for specific agent tasks in Section 5.2. Next, we introduce the birth of SeeClick, including the formalization of GUI grounding task, the construction of continual pre-training data, and training details.

3.1 GUI GROUNDING FOR LVLMs

As GUI grounding is the core capability of SeeClick, we first elucidate how to train LVLM for language generation to perform grounding tasks. Given an interface screenshot s and a collection of elements $\{(x_i, y_i)|_i\}$ on it, where x_i denotes the textual description of the i -th element and y_i indicates the element’s location (represented as a bounding box or point). As depicted in Figure 2(a), LVLM predicts the location of the element y based on the interface screenshot s and its textual description x , i.e. calculating $p(y|s, x)$.

A potential challenge is how LVLMs predict numerical coordinates in a language generation format. Traditional methods (Chen et al., 2021; Wang et al., 2023; Shaw et al., 2023) divided the image into 1000 bins, and creating a new 1,000-token vocabulary $\{< p0 >, < p1 >, \dots, < p999 >\}$ to represent the x and y coordinates. In this work, we adopt a more intuitive manner used in LVLMs (Chen et al., 2023b; Bai et al., 2023), treating numerical values as natural language without any additional tokenization or pre-/post-processing. For instance, in Figure 2(a), for a smartphone screenshot and the instruction “View the new album of Jony J”, we craft a query prompt: “In the UI, where should I click if I want to <instruction>?”. Subsequently, we normally compute the cross-entropy loss between the model output and the groundtruth “click (0.49, 0.40)” to optimize the LVLM.

3.2 DATA CONSTRUCTION

We train SeeClick using three collections of data: web UI data crawled from the internet, mobile UI data reorganized from public datasets and general vision-language instruction-following data.

Web Data. Web UIs, featuring a variety of layouts and design styles across websites, are ideal for training LVLMs’ universal recognition and grounding capabilities across different GUI contexts. We collected approximately 300k web pages from the latest Common Crawl repository to serve as our training data for web UI. For each webpage s , we collect two types of elements from the HTML code as exemplified in Figure 3: (1) elements that display visible text content; and (2) elements with a special “title” attribute that display descriptive text when hovering. This method ensures that we gather a series of interactable elements y and their corresponding instructions x , while encompassing a wide range of text and icon elements. In addition to the grounding task $p(y|s, x)$, we also include web OCR task $p(x|s, y)$, predicting text description based on coordinates.



Figure 3: Example of two types of elements automatically collected from the webpage.

Mobile Data. For mobile UI, we include three types of data: widget captioning, mobile UI grounding, and mobile UI summarization. The widget captioning dataset provides language descriptions for mobile UI elements; for example, the description “play music” for the play button on a music player interface. We utilize the training split of the dataset provided in (Li et al., 2020b), containing nearly

20k screenshots, 40k widgets and 100k descriptions. We derive mobile UI grounding data by reversing the process of widget captioning, treating language descriptions as instructions and corresponding widgets as target elements. To improve diversity, we also incorporate the automatically collected elements and instructions from RICO (Li et al., 2020a). The mobile data involves diverse elements and instructions, facilitating the generalization of SeeClick’s grounding proficiency to diverse GUI contexts. We finally include mobile UI summarization data (Wang et al., 2021) to enhance overall interface comprehension.

General Data. To maintain LVLM’s general capacities on natural images, we incorporate the general vision-language instruction-following data collected with GPT-4 from LLaVA (Liu et al., 2023a), covering conversation, detailed description, and complex reasoning.

Finally, we mix the data above and craft 30 task-specific prompts for each added GUI task, resulting in a 1M dataset to train SeeClick.

3.3 TRAINING DETAILS

We built SeeClick through continual pre-training on a recent advanced LVLM, Qwen-VL (Bai et al., 2023), which possesses grounding capabilities and a higher resolution of 448*448. We train Qwen-VL on the constructed dataset (Section 3.2) for about 10k steps (around 1 epoch) to obtain our GUI base model SeeClick. During training, we employ LoRA (Hu et al., 2021) to fine-tune both the visual encoder and LLM. Further details and task examples are provided in Appendix B.

4 GUI GROUNDING BENCHMARK: *ScreenSpot*

We recognize GUI grounding proficiency as essential for constructing a universal visual GUI agent. However, the constrained capabilities of earlier vision-language models resulted in limited attention on GUI grounding, with few existing studies (Li et al., 2021; Li & Li, 2022; Zhang et al., 2023) largely confined to an Android dataset (Deka et al., 2017) collected in 2017.

To address this research gap, we introduce *ScreenSpot*, an up-to-date, realistic grounding evaluation benchmark encompassing various GUI platforms. It is designed to assess vision-language models’ ability to locate screen elements based on human instructions (Figure 2(b) provides some examples). *ScreenSpot* has two distinctive features: (1) Various GUI platforms. It includes over 600 interface screenshots from mobile (iOS, Android), desktop (macOS, Windows), and web platforms, along with 1200+ instructions and corresponding actionable elements; (2) Icons/Widgets. *ScreenSpot* includes a substantial number of icon/widget type elements in each GUI, which is more challenging to locate than texts (statistics are in Figure 4). See Appendix C.1 for more examples.

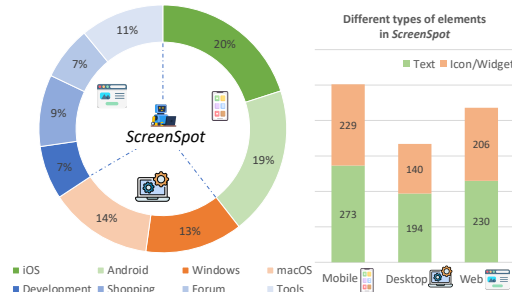


Figure 4: Statistic of our proposed GUI grounding benchmark *ScreenSpot*. The left illustrates the diverse GUI platforms included. The right displays the types of elements within each GUI category.

To evaluate vision-language models’ effectiveness in real-world scenarios, *ScreenSpot* is carefully curated to ensure the samples are novel and not included in existing training resources. We recruited experienced annotators to collect new GUI interfaces and label instructions along with the bounding boxes for actionable elements. For mobile and desktop, annotators were asked to select commonly used apps and operations; for web, we chose several types of websites (development, shopping, forum, and tools) from the recent web environment WebArena (Zhou et al., 2023).

5 EXPERIMENTS

In this section, we first evaluate the GUI grounding capabilities of representative LVLMs and our proposed SeeClick. Subsequently, we adapt SeeClick to mobile and web agent tasks, analyzing the correlation between the advanced grounding capacity and downstream task performance, while exploring the potential of purely vision-based GUI agents.

LVLMs	Model Size	GUI Specific	Mobile		Desktop		Web		Average
			Text	Icon/Widget	Text	Icon/Widget	Text	Icon/Widget	
MiniGPT-v2	7B	✗	8.4%	6.6%	6.2%	2.9%	6.5%	3.4%	5.7%
Qwen-VL	9.6B	✗	9.5%	4.8%	5.7%	5.0%	3.5%	2.4%	5.2%
GPT-4V	-	✗	22.6%	24.5%	20.2%	11.8%	9.2%	8.8%	16.2%
Fuyu	8B	✓	41.0%	1.3%	33.0%	3.6%	33.9%	4.4%	19.5%
CogAgent	18B	✓	67.0%	24.0%	74.2%	20.0%	70.4%	28.6%	47.4%
SeeClick	9.6B	✓	78.0%	52.0%	72.2%	30.0%	55.7%	32.5%	53.4%

Table 1: Results of different LVLMs on *ScreenSpot*. The best results in each column are **bold**. Benefit from efficient GUI grounding pre-training, SeeClick significantly enhanced LVLMs’ ability to locate GUI elements following instructions, and surpassed CogAgent with a smaller model size.

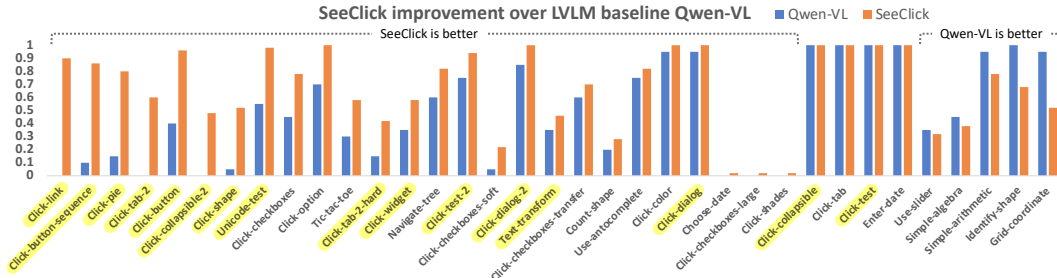


Figure 5: Comparison of SeeClick and LVM baseline Qwen-VL on MiniWob. Tasks with yellow shadows indicate those with dynamic webpage layouts and element positions (closer to the real-world application of GUI agents, see more details in appendix Figure 11). SeeClick outperformed Qwen-VL in most tasks, especially in those with dynamic elements, highlighting the effectiveness of our GUI grounding pre-training.

5.1 GUI GROUNDING ON *ScreenSpot*

As the foundation of visual GUI agents, GUI grounding has not received adequate attention in current LVLMs evaluations (Liu et al., 2023b; Yu et al., 2023). Therefore, we compare various LVLMs on our specially constructed GUI grounding benchmark *ScreenSpot*.

Compared LVLMs & Evaluation. We primarily evaluated two types of LVLMs: (1) Generalist LVLMs capable of tasks such as dialogue, recognition and grounding, including MiniGPT-v2 (Chen et al., 2023a), Qwen-VL (Bai et al., 2023) and GPT-4V; (2) Recently introduced LVLMs specifically designed for GUI tasks, including Fuyu (Bavishi et al., 2023) and CogAgent (Hong et al., 2023).

Considering that GUI agents require clicking on the correct position, we calculate the click accuracy as the metric, defined as the proportion of test samples where the model’ predicted location falls in the groundtruth element bounding box (Li et al., 2022; Zhang et al., 2023).

Results. As shown in Table 1, while generalist LVLMs have excelled in natural image grounding, their GUI grounding performance on *ScreenSpot* is poor due to the significant differences between GUIs and natural images. Even GPT-4V struggles with accurately locating screen elements.

In comparison, LVLMs trained specifically on GUI data have significant improvements. SeeClick achieved the best average performances across three GUI platforms and two types of elements, even with fewer model parameters than CogAgent. This demonstrates the efficiency of our GUI grounding pre-training; with the rich UI elements and diverse instructions collected from the web and mobile, SeeClick quickly learns to understand human instructions for element localization, even in completely unseen scenarios like iOS and desktop. SeeClick exhibits slightly inferior performance in locating text within desktop and web compared to CogAgent, possibly due to lower resolution and much smaller training data. Notably, all models struggle with locating icons/widgets, highlighting the difficulty of identifying and grounding non-text elements on GUIs, which is the unique challenge posed by *ScreenSpot*. More details about evaluation on *ScreenSpot* is in Appendix C.

5.2 VISUAL GUI AGENT TASKS

This section explores the capabilities of SeeClick when applied to three downstream GUI agent tasks: MiniWob, AITW, and Mind2Web. Across these tasks, with provided instructions and memory of previous actions, SeeClick determines the next action solely by observing interface screenshots. The detailed task settings, action formats and interaction examples for SeeClick are in Appendix D.

5.2.1 MINIWOB

MiniWob Shi et al. (2017) comprises about 100 types of web automation tasks, where the agent is asked to interact with a simplified web environment to accomplish human instructions. Existing open-source training data often lacks corresponding interface screenshots (Furuta et al., 2023). Therefore, we rollout 50 successful episodes using an LLM strategy for each task in (Zheng et al., 2023), resulting in a 2.8K episodes dataset to train SeeClick.

Compared Methods & Evaluation. We compared SeeClick with a range of offline training methods. Among these, the state-of-the-art method WebGUM (Furuta et al., 2023) uses screenshots as auxiliary information but still selects HTML elements as actions. Pix2Act (Shaw et al., 2023) is the only prior vision-based approach, trained with extensive demonstration data to perform actions. To verify the effectiveness of SeeClick, we also report the results of the LVLM baseline Qwen-VL when trained with the same 2.8K dataset.

Due to the variance in evaluation task sets among different methods (Liu et al., 2018; Furuta et al., 2023; Shaw et al., 2023), for fairness, we report performance in two groups based on the overlapping MiniWob tasks. We compute the success rate over 50 random seeds for each task and then compute the mean over all tasks as the final score. We provided task-wise scores in Appendix D.2.

Results. As depicted in Table 2, purely vision-based SeeClick surpassed strong baselines with substantially less training data. Notably, with an equivalent amount of 2.8K training data, it outperformed the offline sota WebGUM, which uses both HTML and screenshots as input. Moreover, thanks to LVLM’s powerful reasoning and planning abilities and our GUI grounding pre-training, SeeClick exceeded the sota visual method Pix2Act, using less than 0.3% training data.

Furthermore, SeeClick significantly surpassed the LVLM baseline Qwen-VL by nearly 20 percentage points, underscoring the importance of GUI grounding in boosting LVLM’s performance. To analyze in detail, we provide task-level comparisons between SeeClick and Qwen-VL on 35 MiniWob tasks in Figure 5. SeeClick notably excelled in tasks with dynamic interface layouts and element positions, confirming our hypothesis that general LVLMs struggle with accurately clicking, and SeeClick markedly improves this aspect.

5.2.2 AITW

We evaluate SeeClick in smartphone environments with a recently proposed Android automation dataset Android In The Wild (AITW), which encompasses 30k distinct instructions and corresponding 715k human-annotated operation trajectories. Previous approaches (Rawles et al., 2023) split train/val/test in episode-wise. However, this setting poses a clear risk of overfitting due to: (1) instructions in the test set have appeared during training, and (2) an average of 20 similar trajectories per instruction. In this work, we propose to split the dataset instruction-wise, which avoids overfitting and reflects the performance of agents on unseen instructions. Specifically, we selected

Methods	Modality	Dataset	Score
Compared with text-based models over 45 tasks			
CC-Net (SL)	DOM+Image	2.4M	35.6%
WebN-T5	HTML	12K	55.2%
MM-WebN-T5	HTML+Image	347K	63.4%
WebGUM	HTML+Image	2.8K	65.5%
WebGUM	HTML+Image	347K	86.1%
SeeClick	Image	2.8K	<u>73.6%</u>
Compared with vision-based models over 35 tasks			
CC-Net (SL)	Image	2.4M	23.4%
Pix2Act	Image	1.3M	64.6%
Qwen-VL	Image	2.8K	48.4%
SeeClick	Image	2.8K	<u>67.0%</u>

Table 2: Average scores of different methods on MiniWob. The best results in each setting are **bold**. Methods achieving the highest performance with limited data are underlined. SeeClick outperforms a range of offline training methods as a purely vision-based model.

Methods	Modality	General	Install	GoogleApps	Single	WebShopping	Overall	ClickAcc
ChatGPT-CoT	Text	5.9	4.4	10.5	9.4	8.4	7.7	-
PaLM2-CoT	Text	-	-	-	-	-	39.6	-
GPT-4V	Image	41.7	42.6	49.8	72.8	45.7	50.5	-
Qwen-VL	Image	49.5	59.9	46.9	64.7	50.7	54.3	57.4
SeeClick	Image	54.0	66.4	54.9	63.5	57.6	59.3	66.4

Table 3: Average scores on AITW. ClickAcc calculates the accuracy of click operation. The best results in each column are **bold**. SeeClick exhibits the best performance among competing baselines.

Methods	w/o HTML	Cross-Task			Cross-Website			Cross-Domain		
		Ele.Acc	Op.F1	Step SR	Ele.Acc	Op.F1	Step SR	Ele.Acc	Op.F1	Step SR
MindAct (gen)	✗	20.2	52.0	17.5	13.9	44.7	11.0	14.2	44.7	11.9
MindAct	✗	55.1	75.7	52.0	42.0	65.2	38.9	42.1	66.5	39.6
GPT-3.5	✗	20.3	56.6	17.4	19.3	48.8	16.2	21.6	52.8	18.6
GPT-4	✗	41.6	60.6	36.2	35.8	51.1	30.1	37.1	46.5	26.4
Qwen-VL	✓	15.9	86.7	13.3	13.2	83.5	9.2	14.1	84.3	12.0
SeeClick	✓	<u>28.3</u>	87.0	<u>25.5</u>	<u>21.4</u>	80.6	<u>16.4</u>	<u>23.2</u>	84.8	<u>20.8</u>

Table 4: Comparison of different methods on Mind2Web. The best results in each column are **bold**. Improvements of SeeClick over the LVLm baseline are underline. The GUI grounding pre-training nearly doubled the step success rate of SeeClick on real websites.

545/688/306/700/700 instructions from General/Install/GoogleApps/Single/WebShopping respectively, and retained only one annotated trajectory for each instruction. We selected 80% for training and the remaining for testing in each subset. The result of original split and more details are in Appendix D.3.

Compared Methods & Evaluation. We compare SeeClick with two types of baselines: (1) api-based LLMs such as ChatGPT-CoT (Zhan & Zhang, 2023), PaLM2-CoT (Rawles et al., 2023) and the latest GPT-4V (Yan et al., 2023); (2) our trained LVLm baseline Qwen-VL.

We follow Rawles et al. (2023) to adopt the screen-wise action matching score as the main metric and additionally compute the click accuracy (ClickAcc), which calculates the accuracy when both reference and prediction are click operations.

Results. As illustrated in Table 3, SeeClick achieved the best average performance among both API-based LLMs and trained LVLms. Specifically, SeeClick exhibited a 9% increase in click accuracy over Qwen-VL, supporting the idea that GUI grounding enhances agent task performance through precise clicking.

5.2.3 MIND2WEB

To assess SeeClick’s capabilities in websites, we utilize the recently introduced Mini2Web dataset (Deng et al., 2023). Mind2Web comprises over 2000 open-ended tasks collected from 137 real websites, each with a high-level instruction and a corresponding human action trajectory. Mind2Web was originally designed to evaluate text-based web agents, which select actionable elements from simplified HTML in each step. This work explores visual web agents that predict click positions directly from screenshots. For this purpose, we parsed screenshots and target element bounding boxes from the raw dump of Mind2Web for training and testing. To the best of our knowledge, this is the first attempt of web agents relying solely on screenshots as inputs for navigating real websites.

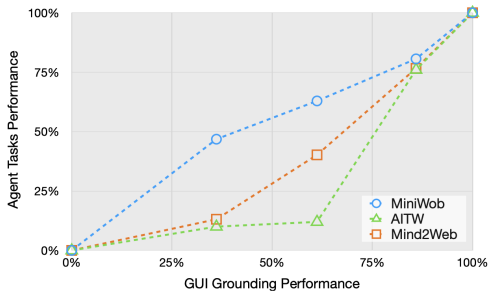
Compared Methods & Evaluation. We compare with html-based web agents Mind2Act (Deng et al., 2023) and our visual baseline Qwen-VL. Mind2Act employs a two-stage method, where a small LM first generates candidate elements from raw HTML, then a large LM selects the target via multi-choice QA; Mind2Act (gen) directly generates the target element instead. GPT-3.5 and GPT-4 adopt the same formulation of Mind2Act and include three demonstrations for in-context learning.

We calculate element accuracy (Ele.Acc), Operation F1 (Op.F1) and step success rate (Step SR). For vision-based methods, a prediction is considered correct if the predicted coordinate falls in the target element’s bounding box. All other settings are following (Deng et al., 2023).

Results. As displayed in Table 4, SeeClick nearly doubled the Ele.Acc and Step SR compared to Qwen-VL. This indicates that SeeClick’s improvement in GUI grounding directly correlates with enhanced performance in web agent tasks. HTML-based methods yield lower Op.F1 as around 20% of groundtruth elements are filtered out during candidate generation. Although SeeClick can operate without extra HTML information, its performance lags behind sota HTML-based methods, as predicting click coordinates on web pages is much more difficult than choosing from HTML candidates. This highlights the difficulty of grounding in intricate interfaces, suggesting the substantial room for improvement needed for visual agents that can apply to the real world.

5.2.4 GROUNDING AND AGENT PERFORMANCE

To investigate the correlation between grounding and agent performance, we examined average scores of several SeeClick’s checkpoints on *ScreenSpot* and three downstream tasks, as depicted in Figure 6. Enhanced GUI grounding capacity boosts downstream task performance, highlighting its crucial role in developing advanced visual GUI agents.



	MiniWob	AITW	Mind2web
Qwen-VL _{separate}	48.4	54.3	11.5
SeeClick _{separate}	67.0	59.3	20.9
SeeClick _{unified}	64.1	57.1	19.5

Table 5: Separate vs unified training performance.

Figure 6: The correlation between agent tasks performance improvement and enhanced grounding ability.

5.2.5 SEECCLICK AS UNIFIED GUI AGENT

To explore SeeClick’s potential for unified operation across GUIs, such as mobile and desktop, we experimented with joint training on three downstream tasks. As shown in Table 5, the unified model exhibited a slight performance decline, possibly due to the distinct interface of different tasks.

6 CONCLUSION

In this paper, we introduce a visual GUI agent - SeeClick, which only relies on screenshots for GUI task automation. We found a key challenge in developing such visual GUI agents: GUI grounding - the capacity to accurately locate screen elements based on human instructions. To address this challenge, we propose to enhance SeeClick via GUI grounding pre-training, and devise methods to automate the curation of GUI grounding data from web and mobile. For benchmarking the progress in GUI grounding, we created *ScreenSpot*, the first realistic evaluation dataset encompassing mobile, desktop, and web platforms. Results on *ScreenSpot* demonstrate a significant improvement of SeeClick over LVLm baselines. Moreover, comprehensive evaluations across three GUI automation tasks consistently support our finding that advancements in GUI grounding directly correlated with improved performance in downstream agent tasks.

LIMITATIONS

SeeClick currently simplifies the GUI action space to mainly focus on clicking and typing, excluding complex actions like dragging and double-clicking. Additionally, limited by the performance of open-source LVLms, training on agent-specific data is necessary for SeeClick to execute multi-step tasks on interfaces like mobile and computer.

REFERENCES

- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023. URL <https://arxiv.org/pdf/2308.12966>.
- Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani, and Sağnak Taşirlar. Introducing our multimodal models, 2023. URL <https://www.adept.ai/blog/fuyu-8b>.
- Andrea Burns, Deniz Arsan, Sanjna Agrawal, Ranjitha Kumar, Kate Saenko, and Bryan A Plummer. A dataset for interactive vision-language navigation with unknown command feasibility. In *European Conference on Computer Vision*, pp. 312–328. Springer, 2022. URL <https://arxiv.org/pdf/2202.02312>.
- Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. Minigpt-v2: large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*, 2023a. URL <https://arxiv.org/pdf/2310.09478>.
- 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*, 2023b. URL <https://arxiv.org/pdf/2306.15195>.
- Ting Chen, Saurabh Saxena, Lala Li, David J Fleet, and Geoffrey Hinton. Pix2seq: A language modeling framework for object detection. In *International Conference on Learning Representations*, 2021. URL <https://arxiv.org/pdf/2109.10852>.
- Biplab Deka, Zifeng Huang, Chad Franzen, Joshua Hibsichman, Daniel Afegan, Yang Li, Jeffrey Nichols, and Ranjitha Kumar. Rico: A mobile app dataset for building data-driven design applications. In *Proceedings of the 30th annual ACM symposium on user interface software and technology*, pp. 845–854, 2017. URL <https://dl.acm.org/doi/pdf/10.1145/3126594.3126651>.
- Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and Yu Su. Mind2web: Towards a generalist agent for the web. *arXiv preprint arXiv:2306.06070*, 2023. URL <https://arxiv.org/pdf/2306.06070>.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020. URL https://arxiv.org/pdf/2010.11929.pdf?fbclid=IwAR1NafJDhZjkARvCswpV6kS9_hMa0ycvzwhlCb7cqAGwgzComFXcScxgA8o.
- Hiroki Furuta, Ofir Nachum, Kuang-Huei Lee, Yutaka Matsuo, Shixiang Shane Gu, and Izzeddin Gur. Multimodal web navigation with instruction-finetuned foundation models. *arXiv preprint arXiv:2305.11854*, 2023. URL <https://arxiv.org/pdf/2305.11854>.
- Difei Gao, Lei Ji, Zechen Bai, Mingyu Ouyang, Peiran Li, Dongxing Mao, Qinchen Wu, Weichen Zhang, Peiyi Wang, Xiangwu Guo, et al. Assistgui: Task-oriented desktop graphical user interface automation. *arXiv preprint arXiv:2312.13108*, 2023. URL <https://arxiv.org/pdf/2312.13108>.
- Izzeddin Gur, Ulrich Rueckert, Aleksandra Faust, and Dilek Hakkani-Tur. Learning to navigate the web. In *International Conference on Learning Representations*, 2018. URL <https://arxiv.org/pdf/1812.09195>.
- Izzeddin Gur, Hiroki Furuta, Austin Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck, and Aleksandra Faust. A real-world webagent with planning, long context understanding, and program synthesis. *arXiv preprint arXiv:2307.12856*, 2023. URL <https://arxiv.org/pdf/2307.12856>.

- Wenyi Hong, Weihang Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxiao Dong, Ming Ding, et al. Cogagent: A visual language model for gui agents. *arXiv preprint arXiv:2312.08914*, 2023. URL <https://arxiv.org/pdf/2312.08914>.
- Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2021. URL <https://arxiv.org/pdf/2106.09685.pdf%20%20>.
- Geunwoo Kim, Pierre Baldi, and Stephen McAleer. Language models can solve computer tasks. *arXiv preprint arXiv:2303.17491*, 2023. URL <https://arxiv.org/pdf/2303.17491>.
- Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. Otter: A multi-modal model with in-context instruction tuning. *arXiv preprint arXiv:2305.03726*, 2023. URL <https://arxiv.org/pdf/2305.03726>.
- Gang Li and Yang Li. Spotlight: Mobile ui understanding using vision-language models with a focus. In *The Eleventh International Conference on Learning Representations*, 2022. URL <https://arxiv.org/pdf/2209.14927>.
- Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, et al. Grounded language-image pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10965–10975, 2022. URL http://openaccess.thecvf.com/content/CVPR2022/papers/Li_Grounded_Language-Image-Pre-Training_CVPR_2022_paper.pdf.
- Yang Li, Jiacong He, Xin Zhou, Yuan Zhang, and Jason Baldridge. Mapping natural language instructions to mobile ui action sequences. *arXiv preprint arXiv:2005.03776*, 2020a. URL <https://arxiv.org/pdf/2005.03776>.
- Yang Li, Gang Li, Luheng He, Jingjie Zheng, Hong Li, and Zhiwei Guan. Widget captioning: Generating natural language description for mobile user interface elements. *arXiv preprint arXiv:2010.04295*, 2020b. URL <https://arxiv.org/pdf/2010.04295>.
- Yang Li, Gang Li, Xin Zhou, Mostafa Dehghani, and Alexey Gritsenko. Vut: Versatile ui transformer for multi-modal multi-task user interface modeling. *arXiv preprint arXiv:2112.05692*, 2021. URL <https://arxiv.org/pdf/2112.05692>.
- Evan Zheran Liu, Kelvin Guu, Panupong Pasupat, Tianlin Shi, and Percy Liang. Reinforcement learning on web interfaces using workflow-guided exploration. In *International Conference on Learning Representations*, 2018. URL <https://arxiv.org/pdf/1802.08802>.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Neural Information Processing Systems*, 2023a. URL <https://arxiv.org/pdf/2304.08485>.
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? *arXiv preprint arXiv:2307.06281*, 2023b. URL <https://arxiv.org/pdf/2307.06281>.
- OpenAI. GPT-4 technical report, 2023. URL <https://arxiv.org/pdf/2303.08774.pdf>.
- Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023. URL <https://arxiv.org/pdf/2306.14824>.
- Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy Lillicrap. Android in the wild: A large-scale dataset for android device control. *arXiv preprint arXiv:2307.10088*, 2023. URL <https://arxiv.org/pdf/2307.10088>.
- Peter Shaw, Mandar Joshi, James Cohan, Jonathan Berant, Panupong Pasupat, Hexiang Hu, Urvashi Khandelwal, Kenton Lee, and Kristina Toutanova. From pixels to ui actions: Learning to follow instructions via graphical user interfaces. In *Advances in Neural Information Processing Systems*, 2023. URL <https://arxiv.org/abs/2306.00245>.

- Tianlin Shi, Andrej Karpathy, Linxi Fan, Jonathan Hernandez, and Percy Liang. World of bits: An open-domain platform for web-based agents. In *International Conference on Machine Learning*, pp. 3135–3144. PMLR, 2017. URL <https://proceedings.mlr.press/v70/shi17a.html>.
- Qiushi Sun, Zhangyue Yin, Xiang Li, Zhiyong Wu, Xipeng Qiu, and Lingpeng Kong. Corex: Pushing the boundaries of complex reasoning through multi-model collaboration. *arXiv preprint arXiv:2310.00280*, 2023. URL <https://arxiv.org/pdf/2310.00280>.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. URL <https://arxiv.org/pdf/2302.13971>.
- Bryan Wang, Gang Li, Xin Zhou, Zhourong Chen, Tovi Grossman, and Yang Li. Screen2words: Automatic mobile ui summarization with multimodal learning. In *The 34th Annual ACM Symposium on User Interface Software and Technology*, pp. 498–510, 2021. URL <https://dl.acm.org/doi/pdf/10.1145/3472749.3474765>.
- 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. *arXiv preprint arXiv:2305.11175*, 2023. URL <https://arxiv.org/pdf/2305.11175>.
- Fangzhi Xu, Zhiyong Wu, Qiushi Sun, Siyu Ren, Fei Yuan, Shuai Yuan, Qika Lin, Yu Qiao, and Jun Liu. Symbol-llm: Towards foundational symbol-centric interface for large language models. *arXiv preprint arXiv:2311.09278*, 2023. URL <https://arxiv.org/pdf/2311.09278>.
- An Yan, Zhengyuan Yang, Wanrong Zhu, Kevin Lin, Linjie Li, Jianfeng Wang, Jianwei Yang, Yiwu Zhong, Julian McAuley, Jianfeng Gao, et al. Gpt-4v in wonderland: Large multimodal models for zero-shot smartphone gui navigation. *arXiv preprint arXiv:2311.07562*, 2023. URL <https://arxiv.org/pdf/2311.07562>.
- Zhao Yang, Jiaxuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. Appagent: Multimodal agents as smartphone users. *arXiv preprint arXiv:2312.13771*, 2023a. URL <https://arxiv.org/pdf/2312.13108>.
- Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. The dawn of llms: Preliminary explorations with gpt-4v (ision). *arXiv preprint arXiv:2309.17421*, 9(1):1, 2023b. URL <https://www.stableaiprompts.com/wp-content/uploads/2023/10/Chatgpt-Updates.pdf>.
- Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, et al. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*, 2023. URL https://arxiv.org/pdf/2304.14178.pdf?trk=public_post_comment-text.
- Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*, 2023. URL <https://arxiv.org/pdf/2308.02490>.
- Zhuosheng Zhan and Aston Zhang. You only look at screens: Multimodal chain-of-action agents. *arXiv preprint arXiv:2309.11436*, 2023. URL <https://arxiv.org/pdf/2309.11436>.
- Zhizheng Zhang, Wenxuan Xie, Xiaoyi Zhang, and Yan Lu. Reinforced ui instruction grounding: Towards a generic ui task automation api. *arXiv preprint arXiv:2310.04716*, 2023. URL <https://arxiv.org/pdf/2310.04716>.
- Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. Gpt-4v (ision) is a generalist web agent, if grounded. *arXiv preprint arXiv:2401.01614*, 2024. URL <https://arxiv.org/abs/2401.01614>.

Longtao Zheng, Rundong Wang, Xinrun Wang, and Bo An. Synapse: Trajectory-as-exemplar prompting with memory for computer control. In *NeurIPS 2023 Foundation Models for Decision Making Workshop*, 2023. URL <https://ltzheng.github.io/Synapse/static/Synapse.pdf>.

Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Yonatan Bisk, Daniel Fried, Uri Alon, et al. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023. URL <https://arxiv.org/pdf/2307.13854>.

Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. URL <https://arxiv.org/pdf/2304.10592>.

A APPENDIX

B DETAILS OF SEECCLICK PRE-TRAINING

B.1 PRE-TRAINING TASKS

SeeClick employs pre-training tasks as outlined in Table 6. For the grounding task, we incorporate two forms: predicting center point coordinates (text_2_point) and predicting bounding box (text_2_bbox). For the task of generating text for elements (similar to OCR), we also include two categories: predicting text based on center point (point_2_text, widget captioning) coordinates and based on bounding boxes (bbox_2_text). Our preliminary experiments indicated that predicting points was slightly better than bounding boxes, likely due to the variable sizes of UI elements. Consequently, we increased the proportion of data with point localization. Finally, about 1 million samples are used for the continual pre-training of SeeClick.

Domain	Task	Sample Num
Web	text_2_point	271K
	text_2_bbox	54K
	point_2_text	54K
	bbox_2_text	54K
Mobile	text_2_point	274K
	text_2_bbox	56K
	UI summarization	48K
	widget captioning	42K
General	LLaVA	145K
	Total	1M

Table 6: All training data used by SeeClick.

For tasks involving coordinates, positions are represented as either the point (x,y) or the bounding box of (left, top, right, down), where each value is a decimal number in the range [0,1] indicating the ratio of the corresponding position to the width or height of the image. Figure 7 provides some examples of the pre-training data.

B.2 TRAINING CONFIGURATIONS

We employed the aforementioned data for continual pre-training of Qwen-VL-Chat to develop SeeClick. To enhance LVLm’s understanding of GUI images, we unlocked the gradients of its visual encoder and applied LoRA for fine-tuning. We adopt AdamW as the optimizer and use a cosine annealing scheduler with an init learning rate of 3e-5 and a global batch size of 64. All training takes around 24 hours on 8 NVIDIA A100 GPUs.

C *ScreenSpot* EVALUATION

C.1 SAMPLE SHOWCASE

Figure 10 provides more examples of *ScreenSpot*, which contains a variety of common GUI scenarios for mobile, desktop, and web platforms.

C.2 EVALUATION DETAIL

For comparing baselines, we tested the models’ grounding capabilities using their officially recommended approach. For instance, with CogAgent, we randomly selected prompts from the official set provided, such as “What steps do I need to take to <instruction>? (with grounding)”, then the output coordinates (or the centers of bounding boxes) were taken as predicted points. For GPT-4V, we follow Yang et al. (2023b) to enable it to locate screen elements based on instructions. SeeClick’s predictions with points were marginally better than bounding boxes, thus we selected point prediction for final evaluation.

C.3 SEECCLICK CASE STUDY & ERROR ANALYSIS

Figure 8 presents some examples of SeeClick on *ScreenSpot*. SeeClick can comprehend human instructions and accurately locate screen elements. To conduct a detailed analysis of localization performance, we quantified the distances between predicted points and groundtruth (the center

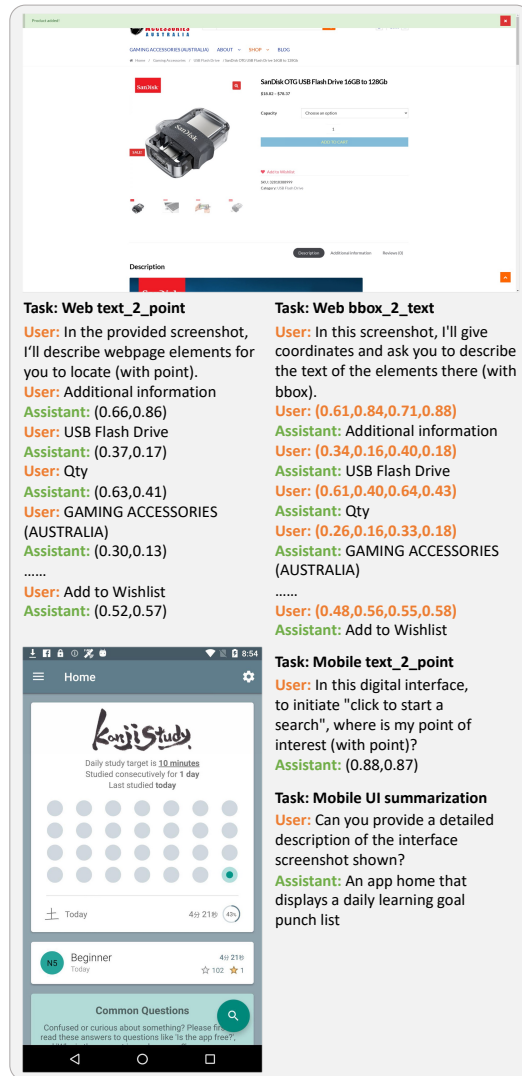


Figure 7: Examples of SeeClick pre-training tasks.

of target elements) in Figure 9. It's noteworthy that even incorrect predictions mostly occur near the target bounding box, suggesting the model recognizes the target but needs improvement in fine-grained localization.

D DOWNSTREAM AGENT TASKS

In this section, we first detail the formulation of SeeClick as a visual GUI agent, then separately introduce the settings for three downstream tasks, and finally show SeeClick's interaction cases with the GUI across these tasks.

D.1 FORMULATION OF SEECCLICK AS VISUAL GUI AGENT

Action Space SeeClick involves common human-UI interaction operations. Following AITW, we assigned an `action_type` id for each action type.

- `click(x, y): 4`. A click action at (x, y) , where each value is a $[0, 1]$ number indicating the ratio of the corresponding position to the width or height of the image.

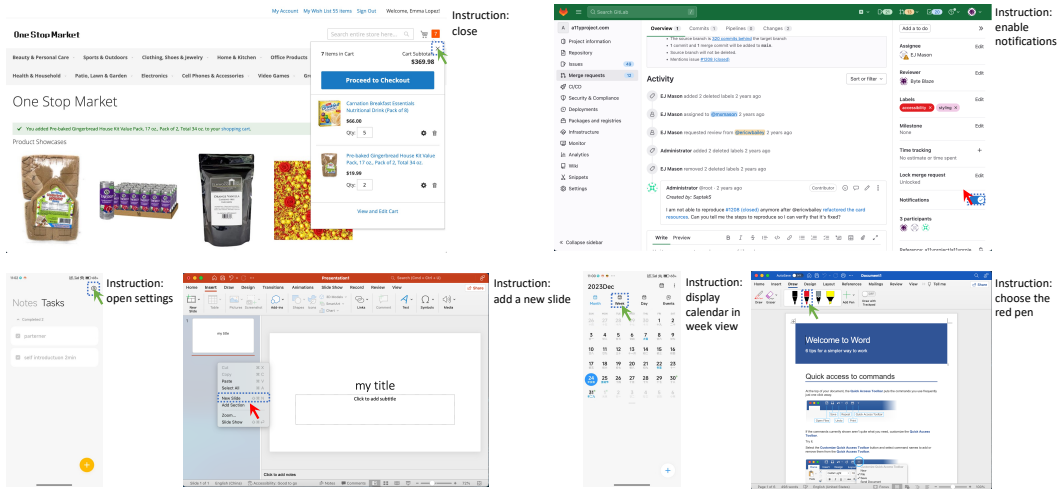


Figure 8: SeeClick on *ScreenSpot*. Blue dashed boxes represent the ground truth bounding boxes, while green and red pointers indicate correct and incorrect predictions.

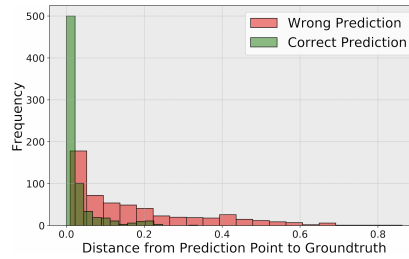


Figure 9: Distance distribution of prediction point to groundtruth. Most incorrect predictions are also close to the answer, suggesting the model recognizes the target but needs improvement in fine-grained localization.

- `type ("typed_text")`: 3. An action of typing a piece of text.
- `select ("value")`: 2. An action for selecting an option from a dropdown menu on a webpage.
- `swipe (direction)`: Swipe actions for the screen, swipe up/down/left/right are assigned the ids 1, 0, 8, and 9 respectively.
- `PRESS BACK`: 5. The action for returning to the previous step.
- `PRESS HOME`: 6. The action for returning to the homepage.
- `PRESS ENTER`: 7. The action of pressing the ENTER key to submit input content.

The first two actions, clicking and typing, are universally applicable across various GUIs. The third action, `select`, is defined according to the specifications in *Mind2Web*. The latter four actions, along with two additional states, `TASK COMPLETE` and `TASK IMPOSSIBLE`, are defined following the *AITW* framework for Android environments.

Agent Formulation SeeClick is an autonomous agent capable of executing human instructions on GUIs. It takes as input the instruction, a screenshot of the current interface and a series of ($k=4$ in our setting) previous actions, to predict the next action to be taken. Specifically, SeeClick uses the following prompt to execute each step of the agent: During training and testing, we organize the data by step into the format described right.

```

<img>Image</img>
User: Please generate the next move according to the
UI screenshot, instruction and previous actions.
Instruction:
<instruction>
Previous actions:
Step1: <step1>
Step2: <step2>
Step3: <step3>
Step4: <step4>
SeeClick: <next action>
    
```


	Gen.	Inst.	GApps.	Sing.	WShop.	Ovr.
Auto-UI	68.2	76.9	71.4	84.6	70.3	74.3
CogAgent	65.4	78.9	75.0	93.5	71.1	76.9
SeeClick	67.6	79.6	75.9	84.6	73.1	76.2

Table 7: Comparison on the origin split of AITW.

D.2 MINIWOB

MiniWob is a classic simplified web agent environment, built on Chrome, allowing low-level operations such as clicking and typing. It comprises around 100 tasks, where each task can template random variants and corresponding instructions controlled by a random seed, creating up to billions of possible task instances. We use 50 successful trajectories for each task provided in (Zheng et al., 2023) for training and test each task with 50 random seeds, following standard practices.

We report the average success rate across random seeds and tasks, automatically provided by the MiniWob environment. A task is considered successfully completed if executed correctly, while incorrect executions or exceeding the maximum number of actions (set as 30 here) are counted as failures. For the baselines in Table 2, we use the task-wise scores provided in their papers to calculate the average for tasks overlapping with SeeClick. We also provided a task-wise comparison in Table 8.

D.3 AITW

AITW is a recently collected dataset for Android smartphone automation, where each sample comprises an instruction and an action trajectory with screenshots. AITW is divided into five subsets: General, Install, GoogleApps, Single, and WebShopping, totally including over 30K instructions and 700K episodes.

Despite AITW’s large scale, as stated in Section 5.2.2, the current train-test split poses a significant risk of overfitting, leading to experimental results that do not accurately reflect an agent’s generalization ability in the real world. We also conducted experiments on SeeClick using the origin split, as shown in Table 7, SeeClick is comparable to CogAgent’s performance. We believe that due to the severe overfitting, designing new agent frameworks or enlarging model size is unlikely to yield much improvements on this split.

To address the aforementioned issue, we propose to divide the train/val/test in an instruction-wise manner. Specifically, we selected 545/688/306/700/700 instructions from the General/Install/GoogleApps/Single/WebShopping subsets (to avoid imbalance in joint training, we randomly chose 700 instructions from Single and WebShopping, which contain a large number of similar instructions), and retained only one annotated episode for each instruction. Next, we allocate 80% for training and the remaining 20% for testing, and select 5*100 episodes to form the validation set from the origin data. The data used for training, validation, and testing will be open-sourced to serve as an effective evaluation.

The other settings are consistent with previous work, calculating a screen-wise matching score that considers both the correctness of the action type and its value (e.g., the click point or typed text), and correlates with the task completion score judged by humans (Rawles et al., 2023).

D.4 MIND2WEB

Mind2Web is a recently proposed dataset for developing generalist web agents for real-world websites, originally designed for text-based agents. Therefore, the origin observation in each step only includes the HTML code of the current webpage. To train and evaluate visual-based agents, we extracted web screenshots and the bounding boxes of target operational elements for each step from Mind2Web’s raw dump. One issue with Mind2Web’s original HTML observation is that it captures the entire page, including scrolling, with its screenshots being long captures (e.g., 1920*12000). Predicting operational positions from such high-resolution long screenshots is impractical for current LVLMS and does not align with human operations. To address this, for target elements not at the top, we

randomly crop around their location, maintaining a consistent screenshot resolution of 1920*1080 for all observed interfaces.

Mind2Web first calculates Element Accuracy (Ele.Acc) which compares the predicted element with groundtruth, and Operation F1 (Op.F1) which calculates the token-level F1 score for the predicted operation. Operation F1 is equivalent to the accuracy of click operations but takes into account the correctness of input values for type and select operations. For our vision-based approach, Element Accuracy is computed as the accuracy of predicted click points falling in the groundtruth elements' bounding box. Then, a step is considered successful (Step SR) if both the predicted element and operation are correct.

D.5 CASE STUDY

MiniWob Figure 11(a) illustrates the difference between static and dynamic layout tasks. Static layout tasks have fixed element positions during training and testing, while dynamic layout tasks display varying interfaces and element positions with instructions, further challenging the agent's ability to accurately locate the target. Figure 11(b) provides examples of SeeClick's interaction with MiniWob. SeeClick relies solely on the interface screenshot for arithmetic, reasoning, etc.

AITW Figure 12 provides SeeClick's operations on AITW. Predictions marked in red below indicate that they were computed as incorrect in AITW. However, some errors occur because the current step's answer is not unique. For example in step 5, the model's predicted input "DuckDuckGo Privacy Browser" is also a potentially correct action.

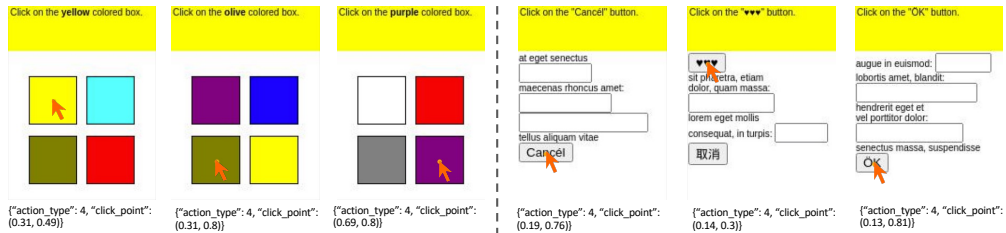
Mind2Web Figure 13 shows several examples of SeeClick on the real-world website benchmark Mind2Web. SeeClick can comprehend instructions and click on the correct elements within complex interfaces.

	CC-Net (SL)	WebN-T5	WebGUM	Pix2Act	Qwen-VL	SeeClick
Choose-date	0.12	0.00	0.13	0.06	0.0	0.02
Click-button	0.78	1.0	1.0	0.32	0.42	0.96
Click-button-sequence	0.47	1.0	1.0	1.0	0.08	0.86
Click-checkboxes	0.32	0.96	1.0	0.99	0.44	0.78
Click-checkboxes-large	0.0	0.22	0.99	1.0	0.0	0.02
Click-checkboxes-soft	0.04	0.54	0.98	0.91	0.06	0.22
Click-checkboxes-transfer	0.36	0.63	0.99	0.76	0.60	0.70
Click-collapsible-2	0.17	0.00	0.95	0.31	0.0	0.48
Click-collapsible	0.81	0.00	0.98	0.80	1.0	1.0
Click-color	0.82	0.27	0.34	0.88	0.96	1.0
Click-dialog	0.95	1.0	1.0	0.12	0.96	1.0
Click-dialog-2	0.88	0.24	0.43	0.73	0.84	1.0
Click-link	0.59	1.0	1.0	0.86	0.0	0.90
Click-option	0.21	0.37	1.0	0.0	0.70	1.0
Click-pie	0.15	0.51	0.99	0.81	0.16	0.80
Click-shades	0.04	0.0	0.0	0.76	0.0	0.02
Click-shape	0.11	0.53	0.72	0.19	0.04	0.52
Click-tab	0.95	0.74	1.0	0.54	1.0	1.0
Click-tab-2	0.27	0.18	0.95	0.52	0.0	0.60
Click-tab-2-hard	0.19	0.12	0.95	0.0	0.16	0.42
Click-test	1.0	1.0	1.0	1.0	1.0	1.0
Click-test-2	0.95	1.0	1.0	1.0	0.72	0.94
Click-widget	0.56	1.0	1.0	0.87	0.38	0.58
Count-shape	0.21	0.41	0.68	0.0	0.20	0.28
Copy-paste	0.04	-	-	-	0.96	0.80
Copy-paste-2	0.01	-	-	-	0.96	0.80
Email-inbox	0.09	0.38	0.99	-	0.08	0.80
Email-inbox-forward-nl	0.0	0.6	1.0	-	0.24	0.74
Email-inbox-forward-nl-turk	0.0	0.33	1.0	-	0.16	0.56
Email-inbox-nl-turk	0.05	0.23	0.98	-	0.40	0.68
Enter-date	0.02	0.0	1.0	0.59	1.0	1.0
Enter-password	0.02	0.97	1.0	-	1.0	1.0
Enter-text	0.35	0.89	1.0	-	1.0	1.0
Enter-text-dynamic	0.39	0.98	1.0	-	0.96	1.0
Focus-text	0.99	1.0	1.0	-	1.0	1.0
Focus-text-2	0.96	1.0	1.0	-	0.84	0.96
Find-word	0.05	-	-	-	1.0	0.10
Grid-coordinate	0.66	0.49	1.0	0.97	0.96	0.52
Guess-number	0.21	0.0	0.11	-	1.0	1.0
Login-user	0.0	0.82	1.0	-	1.0	1.0
Login-user-popup	0.02	0.72	0.99	-	0.86	0.98
Multi-layouts	0.00	0.83	1.0	-	0.44	0.72
Multi-orderings	0.0	0.88	1.0	-	0.42	0.86
Identify-shape	0.68	-	-	0.94	1.0	0.68
Navigate-tree	0.32	0.91	1.0	0.07	0.60	0.82
Search-engine	0.15	0.34	0.96	-	0.56	0.84
Simple-algebra	0.03	-	-	0.99	0.48	0.38
Simple-arithmetic	0.38	-	-	0.67	0.92	0.78
Text-transform	0.19	-	-	0.91	0.36	0.46
Tic-tac-toe	0.32	0.48	0.56	0.76	0.30	0.58
Unicode-test	0.86	-	-	0.64	0.54	0.98
Use-autocomplete	0.07	0.22	0.98	0.95	0.72	0.82
Use-slider	0.18	-	-	0.69	0.38	0.32
Use-spinner	0.47	0.07	0.11	-	0.24	0.16
Read-table	0.01	-	-	-	0.90	0.72
Average	0.336 (55)	0.552 (45)	0.861 (45)	0.646 (35)	0.564 (55)	0.712 (55)

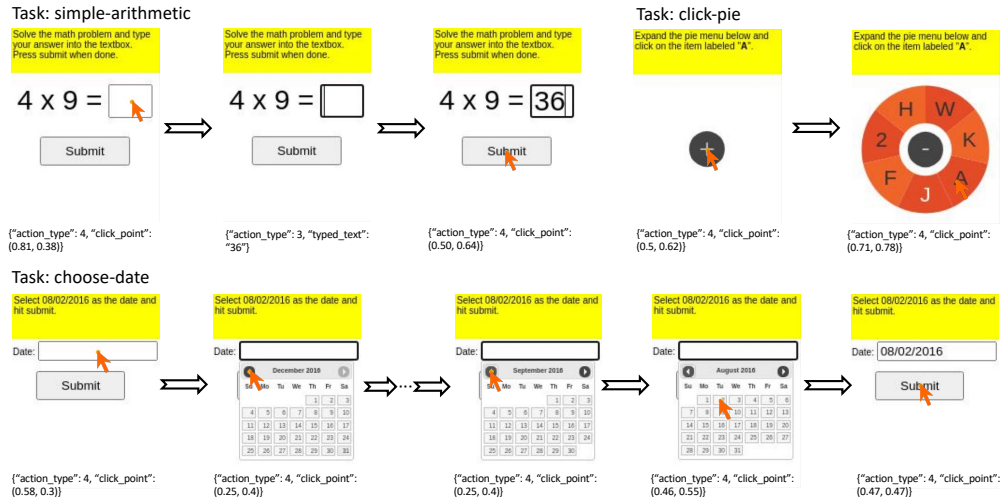
Table 8: Mean scores across 55 MiniWeb tasks.



Figure 10: More examples of GUI grounding benchmark ScreenSpot.



(a) Comparison between static layout (left, click-color) and dynamic layout (right, unicode-test).



(b) Example episodes of SeeClick on MiniWeb tasks.

Figure 11: Example episodes of SeeClick on MiniWeb. The model's prediction output is below the screenshot, with action_type 4 indicating a click and action_type 3 indicating typing.

Instruction: open app "DuckDuckGo Privacy Browser" (install if not already installed) and enter user name: "cleaving@outlook.com" and password: "freighters"

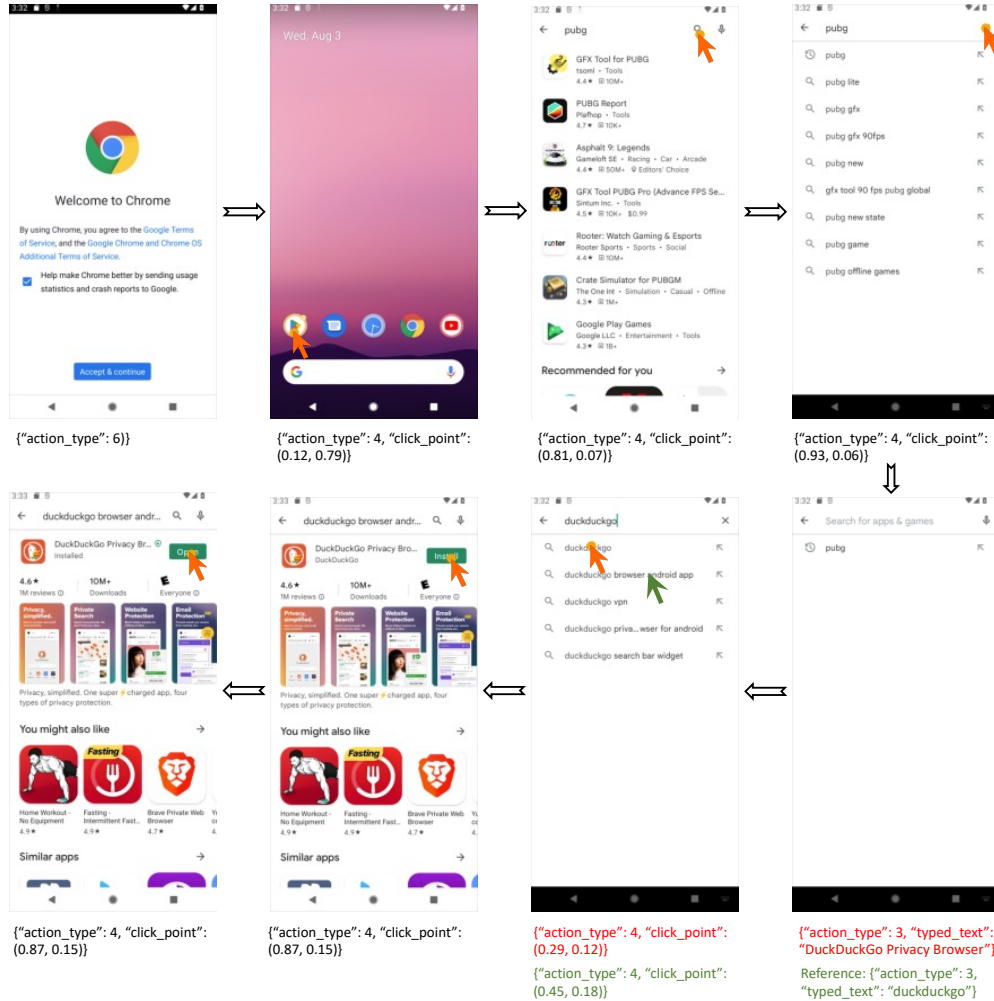


Figure 12: Example episodes of SeeClick on AITW. The model's prediction output is below the screenshot, with action_type 4 indicating a click, action_type 3 indicating typing and action_type 6 indicating PRESS HOME. Steps with the red prediction and green reference indicate a failed step.

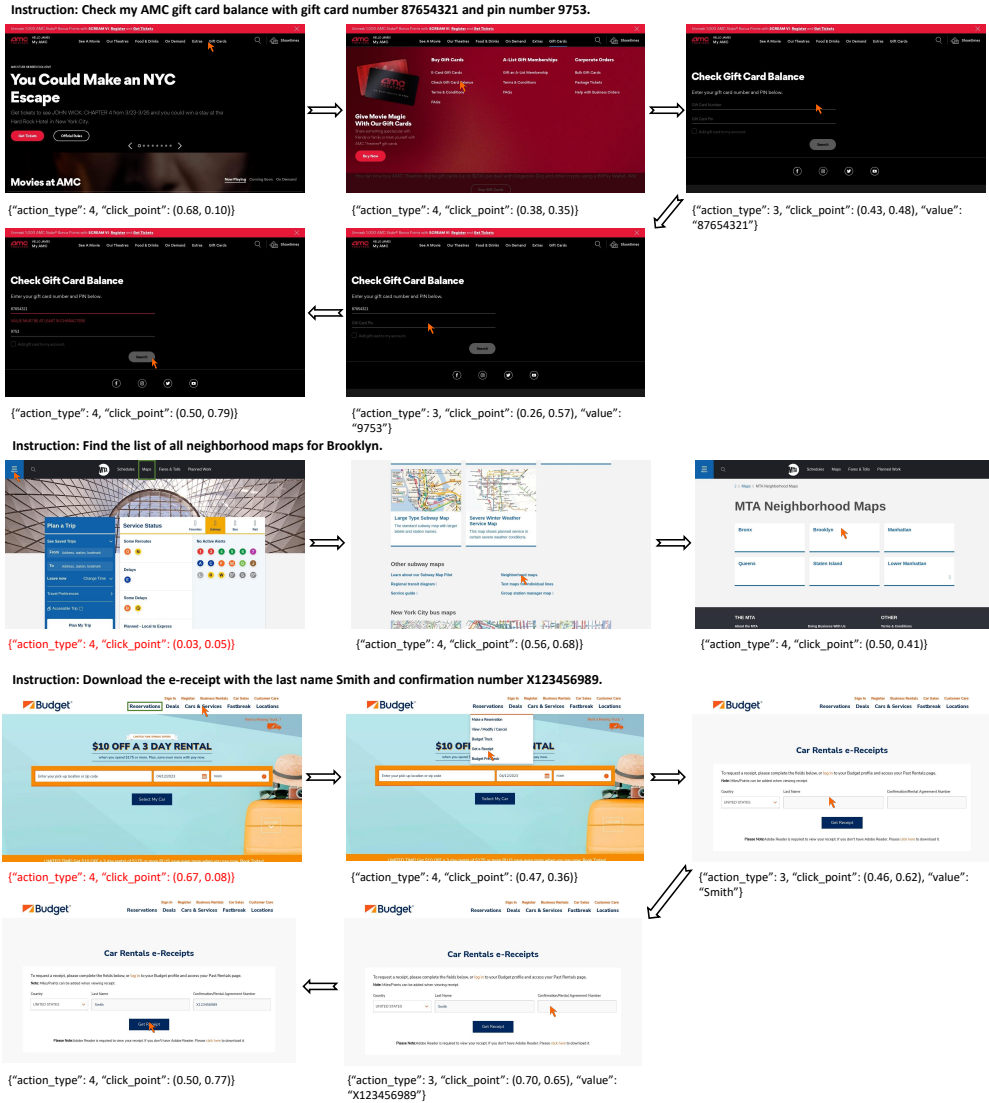


Figure 13: Example episodes of SeeClick on Mind2Web. The model’s prediction output is below the screenshot, with action_type 4 indicating a click and action_type 3 indicating typing. Steps with the red prediction and green reference bounding box indicate a failed step.