

# Explorer: Scaling Exploration-driven Web Trajectory Synthesis for Multimodal Web Agents

Anonymous ACL submission

## Abstract

Recent success in large multimodal models (LMMs) has sparked promising applications of agents capable of autonomously completing complex web tasks. While open-source LMM agents have made significant advances in offline evaluation benchmarks, their performance still falls substantially short of human-level capabilities in more realistic online settings. A key bottleneck is the lack of diverse and large-scale trajectory-level datasets across various domains, which are expensive to collect. In this paper, we address this challenge by developing a scalable recipe to synthesize the largest and most diverse trajectory-level dataset to date, containing over 94K successful multimodal web trajectories, spanning 49K unique URLs, 720K screenshots, and 33M web elements. In particular, we leverage extensive web exploration and refinement to obtain diverse task intents. The average cost is 28 cents per successful trajectory, making it affordable to a wide range of users in the community. Leveraging this dataset, we train **Explorer**, a multimodal web agent, and demonstrate strong performance on both offline and online web agent benchmarks such as Mind2Web-Live, Multimodal-Mind2Web, and MiniWob++. Additionally, our experiments highlight data scaling as a key driver for improving web agent capabilities. We hope this study makes state-of-the-art LMM-based agent research at a larger scale more accessible.

## 1 Introduction

Graphical User Interfaces (GUIs) serve as the primary medium for user interaction across digital environments. Within the GUI environment, LLM-based agents (Su et al., 2024) have shown great potential in automating complex workflows for human users. These agents are designed to operate across diverse interfaces, including the web (Deng et al., 2023; Zhou et al., 2024; Zheng et al., 2024), desktop (Xie et al., 2024; Wu et al., 2024), and

mobile platforms (Rawles et al., 2023; Yan et al., 2023). Navigating modern GUI interfaces, which integrate textual, graphical, and interactive components, typically requires agents to possess visual grounding, long-term planning, and memory management capabilities.

Recent work (Cheng et al., 2024; Gou et al., 2024) has demonstrated the effectiveness of synthetic data for enhancing visual grounding (Gou et al., 2024; Chen et al., 2024a; Kapoor et al., 2024; Chen et al., 2024b) and planning (Xu et al., 2024c; Zhang et al., 2024). Developing end-to-end GUI agents with long-term planning and grounding capabilities requires training on multi-step trajectory data (Xu et al., 2024a,c; Qin et al., 2025). However, existing trajectory datasets are primarily human-annotated (Deng et al., 2023; Li et al., 2024; Lu et al., 2024) or leverage synthetic data just for task proposal curation (Lai et al., 2024; Chen et al., 2024a). And human annotation is expensive to scale for collecting large and diverse training datasets. Therefore, synthetic data has emerged as a promising alternative to human-annotated data (Hartvigsen et al., 2022; Sahu et al., 2022; Ye et al., 2022; Tang et al., 2023; Mukherjee et al., 2023; Mitra et al., 2024). Collecting trajectory-level datasets presents unique challenges: 1) curating a diverse set of task intents at scale, 2) deploying an agent capable of interacting with a real-world environment to complete these tasks through a series of actions, and 3) verifying whether the task is accomplished by the executed action sequence.

Data diversity is essential for equipping generalist web agents with a broad range of skills. Existing work on synthetic web trajectory generation employs self-instruct for task proposal generation (He et al., 2024b). It formulates task proposals from homepages or parametric LLM knowledge, overlooking the richer content available in deeper web pages, which is essential for achieving broader task diversity. Another line of work leverages web

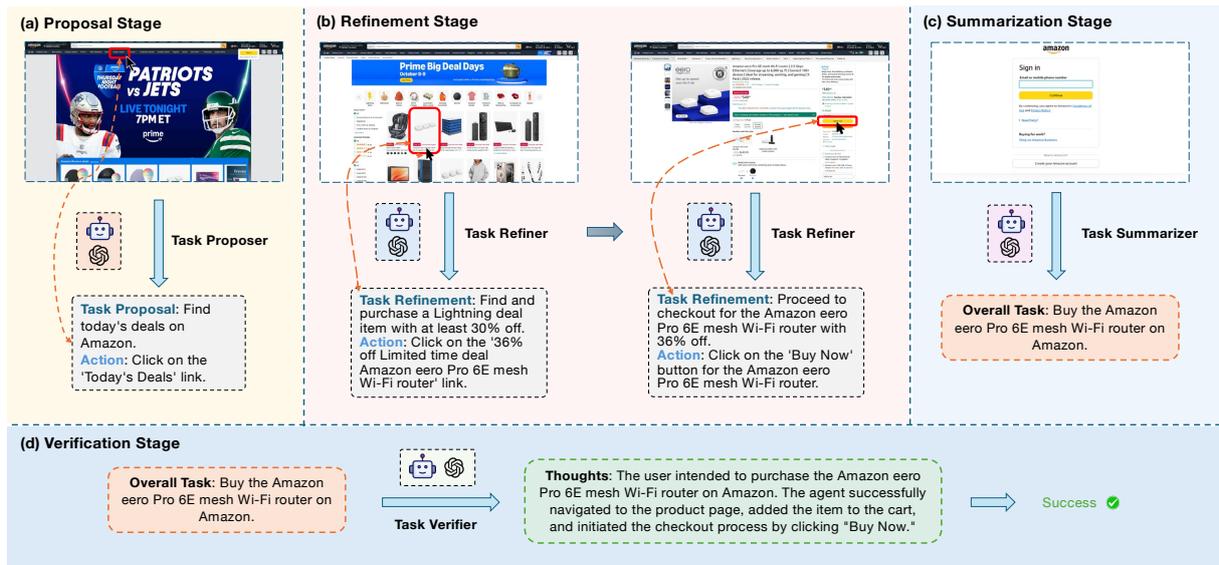


Figure 1: Data Generation Pipeline. The task proposer agent generates an initial task and the first action based on the website homepage. The task is then iteratively refined in subsequent steps by the refiner agent. Finally, the task summarizer agent constructs an overall task description from the action sequence, followed by task verification to assess correctness.

tutorials as a form of supervision for generating web trajectories (Ou et al., 2024; Xu et al., 2024a). While web tutorials effectively cover common daily user tasks, the resulting trajectory data exhibits limited domain diversity in terms of website and domain coverage (Table 1). Additionally, information-seeking tasks remain underrepresented. Due to these limitations, web agents trained on existing synthetic trajectory datasets have not seen much success in more realistic online evaluation settings. To enhance web agents' performance in real-world settings, it is essential to incorporate greater diversity in their training trajectories.

In this work, we develop a *scalable* and *diverse* web trajectory data synthesis recipe for training GUI agent models. Inspired by how humans learn to use the internet, we leverage exploration as a key mechanism for achieving diversity in task intents. We introduce **ExplorER** (**EXPLOR**ation-driven web traject**OR**y g**EN**erato**R**), a framework for systematic web exploration to generate diverse, high-quality trajectory datasets. Unlike prior work that relies on static task proposals, ExplorER dynamically explores web environments to curate diverse, real-world tasks. This exploration-based approach ensures broader task coverage and better generalization to real-world scenarios. We instantiate this framework using popular URLs from several sources, such as Tranco (Pochat et al., 2019) and similarweb.com as seeds. Our dataset comprises

94K diverse web trajectories spanning 49K unique URLs, making it the largest web trajectory dataset to date. Each trajectory is richly annotated with artifacts such as screenshots, raw and set-of-mark (Yang et al., 2023) annotated versions, HTML, and the accessibility tree, enabling comprehensive web agent training. To construct this dataset, we develop a multi-agent pipeline that starts with an abstract task proposal and iteratively refines it into a more specific task through web exploration (Figure 1). Unlike previous approaches, our pipeline generates tasks better grounded in real-world websites, improving task relevance and diversity. To demonstrate the effectiveness of our dataset, we train small language models using just the synthetic data and outperform existing web agent baselines by a significant margin. The main contributions of this work are as follows:

- We develop a scalable and easily customizable multi-agent pipeline for web agent trajectory synthesis. This pipeline leverages exploration as a core mechanism to generate diverse trajectory data, ensuring broad domain coverage and skill diversity in the resulting dataset.
- We leverage this pipeline to generate a diverse and high-quality GUI trajectory dataset consisting of **94K trajectories**, spanning **49K unique URLs** with 720K screenshots and 33M web elements, making it the largest web

	# Trajectories	# Websites	Modality
RUSS (Xu et al., 2021)	80	22	HTML
Mind2Web (Deng et al., 2023)	2350	137	HTML + Screenshot
WebLINX (Lu et al., 2024)	969	155	HTML + Screenshot
GUIAct (Chen et al., 2024a)	2482	121	Screenshot
OpenWebVoyager (He et al., 2024b)	1165	48	A11y tree + Screenshot
AgentTrek (Xu et al., 2024a)	10.4K	127	A11y tree + HTML + Screenshot
NNetnav (Murty et al., 2024)	6K	4	A11y tree + Screenshot
<b>Explorer</b>	<b>94K</b>	<b>49K</b>	<b>A11y tree + Screenshot (raw + SoM) + HTML</b>

Table 1: Comparison to existing web agent benchmarks.

trajectory dataset of this scale.

- We demonstrate the effectiveness of our dataset by training small language models, which achieve strong performance on both online and offline benchmarks, significantly surpassing existing web agent baselines, including those with larger parameter counts.

## 2 Related Work

Recent advances in multimodal language models have facilitated the development of web agents — autonomous systems designed to interact with real-world websites to perform everyday tasks (Deng et al., 2023; Hong et al., 2024; Cheng et al., 2024; Zheng et al., 2024). Early efforts to acquire trajectory data for training web agents primarily relied on crowd-sourcing (Deng et al., 2023; Lu et al., 2024). However, human annotation is cost-prohibitive, prompting the adoption of synthetic data generation approaches to facilitate large-scale data collection. AutoWebGLM (Lai et al., 2024) and GUIAct (Chen et al., 2024a) utilize LLMs to generate task proposals, which human experts subsequently annotate. OpenWebVoyager (He et al., 2024b) employs a web agent to execute auto-generated task descriptions. However, since these task descriptions are generated using LLMs without exploring a website, they fail to capture the full diversity of possible tasks on that website. Another line of work, including Synatra (Ou et al., 2024) and AgentTrek (Xu et al., 2024a), leverages web tutorials to guide web trajectory generation. Meanwhile, concurrent effort (Murty et al., 2024) employs an exploration-based trajectory generation in WebArena’s sandbox, while our work focuses on more realistic web agent evaluation on live websites. To address diversity limitations in prior trajectory synthesis work, we design a bottom-up web trajectory synthesis pipeline that explores websites dynamically while maintaining a coherent high-level task intent. We

direct the reader to Appendix C for a more comprehensive discussion of related work.

## 3 Data Recipe

We design an automatic web trajectory synthesis pipeline that explores websites to generate diverse web trajectories. It utilizes Playwright<sup>1</sup> to execute actions and collect metadata from real-world websites, starting from an initial URL.<sup>2</sup> The metadata includes screenshots, HTML, A11y tree, and actions in grounded and natural language forms. The action space is given in Table B.5 in Appendix.

### 3.1 Website Selection

We use a combination of URL sources to generate the synthetic web trajectories. We obtain the top 100 URLs from [similarweb.com](https://similarweb.com) corresponding to the high-traffic portion of the web with transactional tasks like booking flights, restaurant reservations, government services, sports, entertainment, etc. The Tranco (Pochat et al., 2019) URLs include 49K URLs representing the head portion of the web, which is less trafficked but popular nonetheless. We filter out harmful websites containing violent or explicit content to ensure safety compliance. Overall, we generate 94K trajectories across both sources. The complete data generation takes 50 hours, utilizing 60 parallel processes. The viewport resolution is up to  $1980 \times 1080$ .

### 3.2 Data Generation Pipeline

We aim to develop a generalized pipeline for web exploration to collect diverse web trajectory data. To enhance diversity, we adopt a bottom-up approach, starting with low-level actions and progressively shaping them into high-level task descriptions while maintaining a coherent task intent. In

<sup>1</sup><https://playwright.dev/>

<sup>2</sup>For a 4K subset of trajectories, we instruct GPT-4o to navigate to the target website by formulating a Google search query based on the task description.

Information
View the detailed 7-day weather forecast for Toronto, ON on The Weather Network website. Convert 100 US Dollars to Euros using the XE currency converter. Find directions from Seattle, WA to Bellevue, WA using Bing Maps.
Service
Research the French Bulldog breed on the American Kennel Club website, including its popularity and family life traits. Find the nearest Penske truck rental location in Anaheim, California, and start the reservation process for a truck. Explore and purchase a subscription for the UpToDate Pro Suite on the Wolters Kluwer website.
Entertainment
Find the Basscon presents: Darren Styles EDM event on Eventbrite, save it, and share it on Twitter. View the details of the Photography Competition Winners - Season X and share the article on Twitter.
Shopping
Browse through the fall home decor section on the Target website to explore a variety of fall-themed home decor items. Purchase a three-seat fabric sofa, specifically the UPPLAND Sofa, from IKEA's website.
Travel
Search for flights from Seattle to New York, select travel dates, and explore various flight options. Find the weight of baggage allowance for economy class on qatarairways.

Table 2: Example task descriptions from Explorer.

Metric	Value
# Total trajectories	175K
# Success trajectories	94K
# Unique URLs	49K
Average steps per trajectory	7.7
Average elements per image	46.3
# Tokens	830M
# Elements	33.3M
# Images	720K
Cost per trajectory	\$0.15
Cost per successful trajectory	\$0.28

Table 3: Dataset statistics for Explorer. The number of unique URLs, average steps per trajectory, average elements per image, and number of tokens, elements, and images correspond to the successful trajectories.

the first step, the proposer agent generates an abstract task, which is refined to a more specific task through a refinement process (Figure 1). Since the agents execute actions alongside the refinement process, the generated tasks respect real-world constraints, such as product availability, available color options, and other specifications, ensuring practical applicability. Our pipeline consists of the following LLM-powered agents<sup>3</sup>:

**Task Proposer.** Given a website homepage, including its screenshot and accessibility tree, the task proposer agent generates diverse initial tasks that could be performed on that website. The task descriptions at this stage are instructed to be high-level and abstract versions of the real-world tasks, which will be refined into more specific tasks in later stages. Along with generating the task pro-

<sup>3</sup>We use GPT-4o as the agent backbone throughout the data generation process.

posal, the agent proposes and executes the first action toward completing that task. Furthermore, the agent is instructed to halt upon encountering robot detection such as CAPTCHA verification, login prompts, or payment requests.

**Task Refiner.** The task refiner agent receives the initial task proposal or the refined task description from the previous step, along with the corresponding action history as inputs. It then predicts the next action consistent with the input task description and the updated refined task description while incorporating the complete action history. By iteratively refining the task description after each action, the agent ensures that the updated task remains aligned with the action history.

**Task Summarizer.** This module processes the entire action and screenshot history to predict an overall task description that aligns with the trajectory. The task summary is expected to be high level, *i.e.*, it should describe what the task entails while omitting how it is accomplished.

**Task Verifier.** Inspired by Pan et al. (2024a), the task verifier agent receives the task description and action history, serving as a critic to evaluate whether the trajectory successfully completes the specified task. In addition to the screenshots of the trajectory, it also receives a markdown representation of the last page. This ensures the verifier has the full context of the website's final state, even when the viewport cannot capture all the content. Such automatic evaluation of web trajectories has been widely adopted in prior work (Xu et al., 2024a; He et al., 2024a; Koh et al., 2024). Figure 1 illustrates the above pipeline. The prompts for the

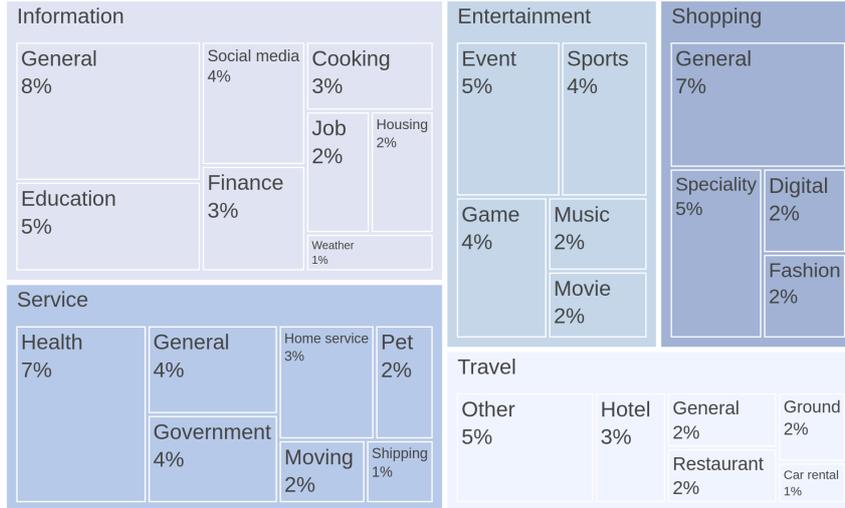


Figure 2: Data composition for Explorer. Its extensive diversity showcases its potential to train end-to-end generalist web agents.

above agents are given in Appendix E.

### 3.3 Dataset Analysis

Explorer comprises web trajectories spanning diverse domains, including services, entertainment, shopping, travel, and information, ensuring broad task diversity. Sample tasks from Explorer are presented in Table 2. Figure 2 visualizes the domain and subdomain distribution, highlighting the dataset’s rich diversity. To the best of our knowledge, Explorer with 94K trajectories is the largest web trajectory dataset of this scale. Table 1 shows a comparison with existing web agent datasets from the literature. The detailed statistics are given in Table 3. Beyond diversity, Explorer is also highly scalable and cost-efficient. Our approach achieves a cost of \$0.28 per successful trajectory, making it approximately  $2\times$  more cost-effective than AgentTrek (Xu et al., 2024a) (which incurs \$0.55 per trajectory) and significantly cheaper than human annotation (Table 4). Unlike human annotation, which requires training crowd workers and continuous quality monitoring, Explorer’s automated generation pipeline eliminates these bottlenecks, ensuring scalability with minimal overhead. By combining diversity, scalability, and cost efficiency, Explorer sets a new benchmark for generating large-scale web trajectory datasets, making it an invaluable resource for training generalist GUI agents.

## 4 Experiments

We use the synthetic trajectories generated by our pipeline to train small multimodal language mod-

Model	Cost per trajectory
Mind2Web (Deng et al., 2023)	\$0.85
AgentTrek (Xu et al., 2024a)	\$0.55
<b>Explorer</b>	<b>\$0.28</b>

Table 4: Cost comparison with other approaches.

els (SLMs) for web agent tasks. To ensure computational efficiency, we select 40K trajectories from the full set for training. We further refine this subset by filtering out trajectories that contain more than two scroll actions to mitigate potential model bias toward excessive scrolling behavior. Finally, we use  $\sim 30$ K trajectories obtained after filtering to fine-tune multimodal language models like Phi-3.5V (Abdin et al., 2024) and Qwen2-VL-7B (Wang et al., 2024a). For brevity, we denote the models trained on Phi-3.5V and Qwen2-VL-7B as Explorer-4B and Explorer-7B, respectively. To test the effectiveness of our data for web-based agentic tasks, we evaluate Explorer-4B and Explorer-7B on Mind2Web-Live (Pan et al., 2024b), Multimodal-Mind2Web (Deng et al., 2023; Zheng et al., 2024), and MiniWob++ (Liu et al., 2018).

**Multimodal-Mind2Web.** Multimodal-Mind2Web is an offline web agent benchmark comprising 2K open-ended tasks spanning 137 websites across 31 domains. Each task comprises a sequence of actions with screenshots, action type, and HTML. We follow the setting in Zheng et al. (2024) and report element accuracy, operation F1, and step success rate (SR) as evaluation metrics.

Model	Avg. Step SR (%)	Completion Rate (%)	Task SR (1) (%)	Full Task SR (%)
<b>API-based Models</b>				
GPT-4o	58.5	52.8	44.6	25.3
GPT-3.5	–	36.5	–	15.4
<b>Open-source Instructed Models</b>				
Mistral-7B-Instr.-0.3 (Jiang et al., 2023)	32.8	29.5	24.1	9.6
Qwen2-72B-Instruct (Bai et al., 2023)	–	<b>40.9</b>	–	15.4
Qwen2-VL-7B (Wang et al., 2024a)	40.2	35.4	<b>34.9</b>	14.5
Phi-3.5V (Abdin et al., 2024)	28.5	23.5	20.5	2.4
<b>Supervised Fine-Tuning</b>				
<b>Explorer-4B</b>	44.0	39.4	31.3	18.1
<b>Explorer-7B</b>	<b>45.3</b>	40.2	<b>34.9</b>	<b>19.3</b>

Table 5: Results on Mind2Web-Live benchmark. Missing values are denoted by –. The results for GPT-4 and Mistral-7B have been reproduced on our Linux servers. The results for GPT-3.5 and Qwen2-72B-Instruct have been taken from Pan et al. (2024b). The full task success rate represents the successful completion of all key nodes for a given task. The average step success rate represents the proportion of completed key nodes, macro-averaged across tasks. The completion rate represents the proportion of completed key nodes, micro-averaged across tasks. Task SR (1) represents task SR with a tolerance of up to one error/key node.

**Mind2Web-Live.** Mind2Web-Live is a benchmark modified from Mind2Web to test web agents on live websites rather than static trajectories. The benchmark evaluates performance using a key-node-based evaluation approach rather than using a golden action sequence, requiring valid trajectories to reach annotated “key nodes” across 104 test tasks in Mind2Web. Since Mind2Web-Live relies on real-world dynamic websites, it encounters robot detection such as reCAPTCHA, which hinders testing (Xu et al., 2024c). To address this, we select a subset of 83 test set tasks that remain consistently accessible throughout our tests. Following Pan et al. (2024b), we report the average step success rate, completion rate, and full task success rate on the test set.

**MiniWeb++.** This benchmark consists of low-level tasks on a single webpage. Typical examples include clicking a sequence of buttons, selecting items from a drop-down list, and filling out a form. We use the subset of 46 tasks used for evaluation in prior work (Zeng et al., 2024; Ou et al., 2024). The final score is obtained by averaging the results of four runs per task. We use the zero-shot evaluation setting, which does not use any environment-specific trajectories for training.

## 5 Results

### 5.1 In-domain Evaluation

As an intrinsic evaluation of the trajectory collection pipeline, we generate 100 test tasks using Explorer, disjoint from the train set. The SLM agents

are tasked with executing the given tasks on live websites while an LLM-as-a-judge verifier (§ 3.2) evaluates the correctness of their actions at the trajectory level. Table 7 shows the results. We observe that the fine-tuned agents significantly outperform their pre-trained counterparts. Thus, using in-domain web trajectory data training helps, which is a valuable sanity check.

### 5.2 Mind2Web-Live Results

We evaluate Explorer-4B and Explorer-7B trained on the synthetic trajectory dataset (Table 5). We make the following observations from the results:

#### **Improvement over base pre-trained models.**

We observe that Explorer-7B yields improvements of 5.1% and 4.8% in average step success rate (SR) and key node completion rate, respectively, compared to the pre-trained Qwen2-VL-7B model. Similarly, Explorer-4B obtains gains of 15.5% and 15.9% in average step SR and key node completion rate, respectively, over its pre-trained counterpart. In terms of full task success rate, Phi-3.5V improves significantly from 2.4% to 18.1%, while Qwen2-VL-7B improves from 14.5% to 19.3%. To the best of our knowledge, this represents the state-of-the-art performance on Mind2Web-Live for models of this size trained exclusively on synthetic data.

**Improvement over higher capacity pre-trained models.** Despite having much fewer parame-

Model	Train Data	Cross-Task			Cross-Website			Cross-Domain			Avg.
		Ele. Acc	Op. F1	Step SR	Ele. Acc	Op. F1	Step SR	Ele. Acc	Op. F1	Step SR	
<b>In-Context Learning</b>											
GPT-3.5		19.4	59.2	16.8	14.9	56.5	14.1	25.2	57.9	24.1	18.3
GPT-4		40.8	63.1	32.3	30.2	61.0	27.0	35.4	61.9	29.7	29.7
SeeAct (Zheng et al., 2024)		46.4	73.4	40.2	38	67.8	32.4	42.4	69.3	36.8	36.5
<b>Supervised Fine-Tuning</b>											
SeeClick-9.6B (Cheng et al., 2024)	Syn. + M2W	26.3	86.2	23.7	21.9	82.9	18.8	22.1	84.1	20.2	20.9
EDGE-9.6B (Chen et al., 2024b)	Syn. + M2W	–	–	30.0	–	–	21.1	–	–	22.4	24.5
MiniCPM-GUI-3.1B (Chen et al., 2024a)	Syn. + M2W	23.8	86.8	20.8	20.3	81.7	17.3	17.9	74.5	14.6	17.6
Scribe-Agent-L-32B (Shen et al., 2024)	Syn. traj.	38.0	52.9	35.6	34.1	52.7	32.5	39.4	54.7	37.3	35.1
AgentTrek-7B (Xu et al., 2024a)	Syn. + M2W	<b>60.8</b>	88.9	<b>55.7</b>	57.6	88.1	51.4	<b>56.0</b>	87.5	52.6	53.2
<b>Explorer-4B</b>	Syn. traj.	36.5	82.9	33.2	44.1	87.7	39.3	42.5	86.3	39.8	37.4
<b>Explorer-4B</b>	M2W	48.1	88.0	44.8	49.1	87.2	45.0	46.9	87.7	44.6	44.8
<b>Explorer-4B</b>	Syn. + M2W	53.4	88.1	50.7	55.6	89.5	51.4	49.8	88.8	47.2	49.8
<b>Explorer-7B</b>	Syn. traj.	43.6	86.6	39.6	48.7	87.7	44.5	47.6	87.2	44.7	43.0
<b>Explorer-7B</b>	M2W	51.8	88.0	48.3	56.3	89.7	52.0	50.9	88.9	48.1	49.5
<b>Explorer-7B</b>	Syn. + M2W	56.5	<b>90.3</b>	53.2	<b>60.5</b>	<b>90.7</b>	<b>56.7</b>	55.7	<b>90.4</b>	<b>53.0</b>	<b>54.3</b>

Table 6: Multimodal-Mind2Web evaluation results. The baseline numbers have been taken from Zheng et al. (2024); Cheng et al. (2024); Chen et al. (2024b,a); Shen et al. (2024). The last column denotes the average step success rates over the three test splits. Explorer significantly outperforms existing GUI agent baselines.

Model	Full Task SR (%)
GPT-4o	16.0
Phi-3.5V	1.0
<b>Explorer-4B</b>	17.0
Qwen2-VL-7B	6.0
<b>Explorer-7B</b>	<b>18.0</b>

Table 7: In-domain evaluation results. The fine-tuned Explorer models achieve significant improvements over their pre-trained counterparts and surpass closed-source LLMs, including GPT-4o.

Model	Accuracy (%)
<i>API-based Models</i>	
GPT-3.5	39.57
GPT-4	53.04
<i>Open-source Instructed Models</i>	
Phi-3.5V	35.87
Qwen2-VL-7B	36.96
Llama3-chat-8B	31.74
Llama3-chat-70B	48.70
<i>Open-source Interactive Data Finetuned Models</i>	
AgentLM-7B (Zeng et al., 2024)	15.65
CodeActAgent-7B (Wang et al., 2024b)	9.78
AgentFlan-7B (Chen et al., 2024c)	20.87
Lemur-chat-70B (Xu et al., 2024b)	21.30
AgentLM-70B (Zeng et al., 2024)	36.52
Synatra-CodeLlama-7B (Ou et al., 2024)	38.20
AgentTrek-7B (Xu et al., 2024a)	45.28
<b>Explorer-4B</b>	46.74
<b>Explorer-7B</b>	<b>53.26</b>

Table 8: Results on MiniWob++ benchmark (Liu et al., 2018) in zero-shot evaluation setting. The baseline numbers correspond to Ou et al. (2024). Explorer outperforms much larger models by a significant margin.

ters, we observe that Explorer-4B outperforms strong baselines such as Mistral-7B-Instruct-0.3 and Qwen2-72B-Instruct in full task SR by margins of 8.5% and 2.7%, respectively. The Phi-3.5V model obtains an 18.1% full task success rate, which is better than GPT-3.5 (15.4%), despite using orders of magnitude fewer parameters. The corresponding results for the entire set of 104 tasks, including unreachable websites, are given in Appendix A. We provide the ablation studies in Appendix A.3 and the error analysis in Appendix A.4.

### 5.3 Multimodal-Mind2Web Results

Following Deng et al. (2023), we obtain the top-50 elements from a pre-trained DeBERTa (He et al., 2021) candidate generation model which are then used to construct the accessibility tree and SoM image inputs. The results are shown in Table 6.

Among baselines, we include API-based models for in-context learning – GPT-3.5, GPT-4, and SeeAct (Zheng et al., 2024). SeeAct is a web agent that performs web tasks using a two-step procedure of action generation and grounding using GPT-4V. Additionally, we include baselines that fine-tune small language models using synthetic data, followed by further fine-tuning on the Mind2Web training set. SeeClick (Cheng et al., 2024) introduces a visual grounding model (Qwen-VL) trained on synthetically-generated grounding data. EDGE (Chen et al., 2024b) synthesizes QA data on webpages to improve the grounded GUI understanding capabilities of MLLMs. ScribeAgent-Large (Shen et al., 2024) and MiniCPM-GUI (Chen

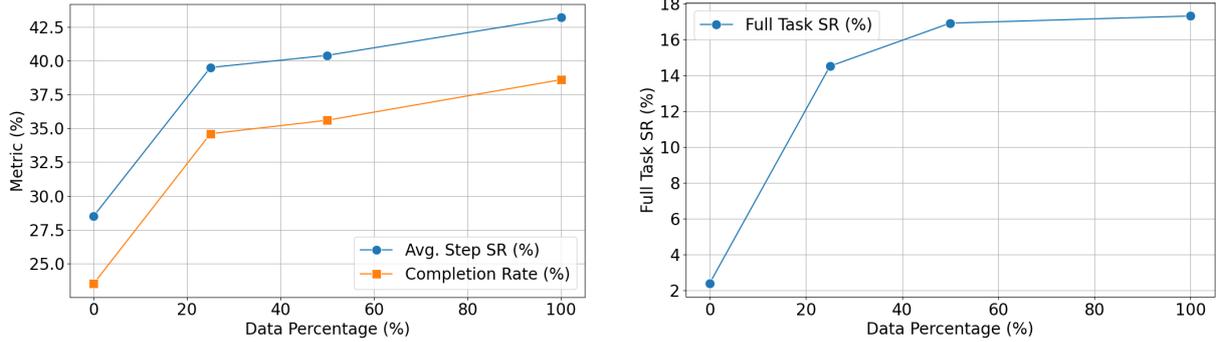


Figure 3: Experiments with data scaling using Explorer-4B on Mind2Web-Live. We experiment with using 100%, 50%, and 25% of the trajectory data. All results are averaged over three runs. All metrics exhibit improvement with increase in data scale.

et al., 2024a) use human-annotated trajectory data to train web agents. AgentTrek (Xu et al., 2024a) is a GUI agent baseline that also utilizes synthetic trajectory data to fine-tune SLMs for Mind2Web, similar to our setting. We observe that Explorer-7B fine-tuned on synthetic data from Explorer plus Mind2Web outperforms all baselines in average step success rate. Notably, it surpasses AgentTrek, which uses the same Qwen2-VL-7B MLLM backbone, highlighting the superior quality of our dataset. The broad domain coverage and task diversity in Explorer contribute to its superior generalization across environments.

#### 5.4 MiniWob++ Results

Table 8 shows the results on the MiniWob++ benchmark in the zero-shot evaluation setting. Among baselines, we have API-based models, in-context learning using open-source LMs, and agentic models like AgentLM (Zeng et al., 2024), CodeActAgent (Wang et al., 2024b), Lemur-Chat (Xu et al., 2024b) and AgentFlan (Chen et al., 2024c) which include web-based demonstrations in their instruction tuning dataset. Synatra-CodeLlama-7B (Ou et al., 2024) and AgentTrek (Xu et al., 2024a) also synthesize web-agent trajectories automatically. We observe that Explorer outperforms GPT-4 and general-purpose agent baselines. Explorer-4B surpasses Synatra-CodeLlama-7B and AgentTrek-7B despite using a much smaller model with 4.2B params, highlighting our synthetic data’s superior quality and potential for out-of-distribution (OOD) generalization.

#### 5.5 Data Scaling Experiments

We conduct experiments with different data scales for Explorer-4B to analyze the impact of training

data size. Specifically, we subsample the original trajectory dataset to utilize 50% and 25% of its original size. Figure 3 presents the resulting performance curves. Our results show that, even with just 25% of the training data, the model exhibits rapid performance gains over the base pre-trained model. Increasing the dataset size further leads to gradual improvements across all reported metrics. However, the increase in the overall task success rate is more gradual compared to the stepwise metrics, as it is a more coarse-grained metric.

## 6 Conclusion

In this work, we introduce Explorer, a scalable framework for synthesizing web trajectories on a large scale. By leveraging thorough web exploration, Explorer ensures diversity in both domains and the skills acquired by web agents. Unlike previous approaches, our framework generates contextually grounded trajectories that adapt to real-world constraints, improving both task relevance and generalization. We instantiate this framework using URLs collected from diverse sources. Explorer outperforms existing web agent baselines by a significant margin on both online and offline web agent benchmarks. Furthermore, our results highlight the critical role of data scale in enhancing web agents’ performance. Future work will focus on extending this framework to encompass a broader range of GUI environments, such as operating systems with diverse applications. GUI agents require specialized skills for different tasks, including information-seeking, operational, and navigation skills. Efficient exploration of the environment to acquire these skills presents another promising avenue for future research.

## 488 Limitations

489 Explorer explores the web environment au-  
490 tonomously, which may occasionally result in inco-  
491 herent tasks. Synthetic data collection using closed-  
492 source LLMs can be costly due to associated API  
493 expenses. While this work serves as a proof of  
494 concept, future research will focus on developing  
495 tailor-made open-source LLMs for this task. Addi-  
496 tionally, some website content remains inaccessible  
497 due to login requirements, leading to insufficient  
498 data for those websites.

## 499 Ethical Considerations

500 The synthetic data collection pipeline proposed in  
501 this paper is intended solely for academic research  
502 on GUI agents, with strict ethical safeguards to pre-  
503 vent unauthorized website interactions. To ensure  
504 ethical compliance and mitigate risks, we prompt  
505 our agents to automatically terminate upon encoun-  
506 tering CAPTCHA verifications, login prompts, or  
507 payment requests, ensuring that no actual transac-  
508 tions or bookings occur. Additionally, we filter  
509 out websites containing violent or explicit content  
510 and strictly adhere to privacy regulations, ensuring  
511 that no personal information is used during action  
512 execution. To enforce responsible data collection,  
513 we monitor a subset of automatically generated tra-  
514 jectories to ensure compliance with website access  
515 policies. Moreover, we distribute the workload  
516 across websites to prevent excessive requests and  
517 minimize the impact on any single domain.

## 518 References

519 Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan,  
520 Jyoti Aneja, Ahmed Awadallah, Hany Awadalla,  
521 Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harki-  
522 rat Behl, et al. 2024. Phi-3 technical report: A highly  
523 capable language model locally on your phone. *arXiv*  
524 *preprint arXiv:2404.14219*.

525 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,  
526 Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei  
527 Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin,  
528 Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu,  
529 Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren,  
530 Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong  
531 Tu, Peng Wang, Shijie Wang, Wei Wang, Sheng-  
532 guang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang,  
533 Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu,  
534 Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingx-  
535 uan Zhang, Yichang Zhang, Zhenru Zhang, Chang  
536 Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang  
537 Zhu. 2023. Qwen technical report. *arXiv preprint*  
538 *arXiv:2309.16609*.

Wentong Chen, Junbo Cui, Jinyi Hu, Yujia Qin, Junjie  
Fang, Yue Zhao, Chongyi Wang, Jun Liu, Guirong  
Chen, Yupeng Huo, et al. 2024a. Guicourse: From  
general vision language models to versatile gui  
agents. *arXiv preprint arXiv:2406.11317*. 539  
540  
541  
542  
543

Xuetian Chen, Hangcheng Li, Jiaqing Liang, Si-  
hang Jiang, and Deqing Yang. 2024b. Edge: En-  
hanced grounded gui understanding with enriched  
multi-granularity synthetic data. *arXiv preprint*  
*arXiv:2410.19461*. 544  
545  
546  
547  
548

Zehui Chen, Kuikun Liu, Qiuchen Wang, Wenwei  
Zhang, Jiangning Liu, Dahua Lin, Kai Chen, and  
Feng Zhao. 2024c. Agent-flan: Designing data and  
methods of effective agent tuning for large language  
models. In *Findings of ACL*. 549  
550  
551  
552  
553

Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu,  
Yantao Li, Jianbing Zhang, and Zhiyong Wu. 2024.  
Seeclck: Harnessing GUI grounding for advanced  
visual GUI agents. In *Proceedings of ACL*. 554  
555  
556  
557

Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen,  
Samual Stevens, Boshi Wang, Huan Sun, and Yu Su.  
2023. Mind2web: Towards a generalist agent for the  
web. In *Proceedings of NeurIPS*. 558  
559  
560  
561

Boyu Gou, Ruohan Wang, Boyuan Zheng, Yanan Xie,  
Cheng Chang, Yiheng Shu, Huan Sun, and Yu Su.  
2024. Navigating the digital world as humans do:  
Universal visual grounding for gui agents. *arXiv*  
*preprint arXiv:2410.05243*. 562  
563  
564  
565  
566

Izzeddin Gur, Hiroki Furuta, Austin V. Huang, Mustafa  
Safdari, Yutaka Matsuo, Douglas Eck, and Aleksan-  
dra Faust. 2024. A real-world webagent with plan-  
ning, long context understanding, and program syn-  
thesis. In *Proceedings of ICLR*. 567  
568  
569  
570  
571

Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi,  
Maarten Sap, Dipankar Ray, and Ece Kamar. 2022.  
Toxigen: A large-scale machine-generated dataset  
for adversarial and implicit hate speech detection. In  
*Proceedings of ACL*. 572  
573  
574  
575  
576

Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu,  
Yong Dai, Hongming Zhang, Zhenzhong Lan, and  
Dong Yu. 2024a. Webvoyager: Building an end-to-  
end web agent with large multimodal models. In  
*Proceedings of ACL*. 577  
578  
579  
580  
581

Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu,  
Hongming Zhang, Tianqing Fang, Zhenzhong Lan,  
and Dong Yu. 2024b. Openwebvoyager: Building  
multimodal web agents via iterative real-world ex-  
ploration, feedback and optimization. *arXiv preprint*  
*arXiv:2410.19609*. 582  
583  
584  
585  
586  
587

Pengcheng He, Xiaodong Liu, Jianfeng Gao, and  
Weizhu Chen. 2021. Deberta: decoding-enhanced  
bert with disentangled attention. In *Proceedings of*  
*ICLR*. 588  
589  
590  
591

592	Wenyi Hong, Weihang Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan Wang, Yuxiao Dong, Ming Ding, and Jie Tang. 2024. <a href="#">Cogagent: A visual language model for GUI agents</a> . In <i>Proceedings of CVPR</i> .	646
593		647
594		648
595		649
596		
597	Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. <i>arXiv preprint arXiv:2310.06825</i> .	650
598		651
599		652
600		653
601		654
602		655
603	Raghav Kapoor, Yash Parag Butala, Melisa Russak, Jing Yu Koh, Kiran Kamble, Waseem AlShikh, and Ruslan Salakhutdinov. 2024. <a href="#">Omniact: A dataset and benchmark for enabling multimodal generalist autonomous agents for desktop and web</a> . In <i>Proceedings of ECCV</i> .	656
604		657
605		658
606		659
607		
608	Jing Yu Koh, Stephen McAleer, Daniel Fried, and Ruslan Salakhutdinov. 2024. Tree search for language model agents. <i>arXiv preprint arXiv:2407.01476</i> .	660
609		661
610		662
611		663
612	Hanyu Lai, Xiao Liu, Iat Long Iong, Shuntian Yao, Yuxuan Chen, Pengbo Shen, Hao Yu, Hanchen Zhang, Xiaohan Zhang, Yuxiao Dong, and Jie Tang. 2024. <a href="#">Autowebglm: A large language model-based web navigating agent</a> . In <i>KDD</i> .	664
613		665
614		
615		
616	Wei Li, William E Bishop, Alice Li, Christopher Rawles, Folawiyo Campbell-Ajala, Divya Tyamagundlu, and Oriana Riva. 2024. <a href="#">On the effects of data scale on UI control agents</a> . In <i>Proceedings of NeurIPS Datasets and Benchmarks Track</i> .	666
617		667
618		668
619		669
620		670
621		
622	Evan Zheran Liu, Kelvin Guu, Panupong Pasupat, Tianlin Shi, and Percy Liang. 2018. <a href="#">Reinforcement learning on web interfaces using workflow-guided exploration</a> . In <i>Proceedings of ICLR</i> .	671
623		672
624		673
625		674
626	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. <a href="#">Visual instruction tuning</a> . In <i>Proceedings of NeurIPS</i> .	675
627		676
628		677
629	Xing Han Lu, Zdenek Kasner, and Siva Reddy. 2024. <a href="#">Weblinx: Real-world website navigation with multi-turn dialogue</a> . In <i>Proceedings of ICML</i> .	678
630		679
631		680
632	Grégoire Mialon, Clémentine Fourrier, Thomas Wolf, Yann LeCun, and Thomas Scialom. 2024. <a href="#">GAIA: a benchmark for general AI assistants</a> . In <i>Proceedings of ICLR</i> .	681
633		682
634		683
635	Arindam Mitra, Luciano Del Corro, Guoqing Zheng, Shweti Mahajan, Dany Rouhana, Andres Codas, Yadong Lu, Wei-ge Chen, Olga Vrousos, Corby Rosset, et al. 2024. <a href="#">Agentinstruct: Toward generative teaching with agentic flows</a> . <i>arXiv preprint arXiv:2407.03502</i> .	684
636		685
637		686
638		687
639		688
640		689
641	Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. 2023. <a href="#">Orca: Progressive learning from complex explanation traces of gpt-4</a> . <i>arXiv preprint arXiv:2306.02707</i> .	690
642		691
643		692
644		693
645		694
		695
		696
		697
		698
	Shikhar Murty, Dzmitry Bahdanau, and Christopher D Manning. 2024. <a href="#">Netscape navigator: Complex demonstrations for web agents without a demonstrator</a> . <i>arXiv preprint arXiv:2410.02907</i> .	646
		647
		648
		649
	Tianyue Ou, Frank F. Xu, Aman Madaan, Jiarui Liu, Robert Lo, Abishek Sridhar, Sudipta Sengupta, Dan Roth, Graham Neubig, and Shuyan Zhou. 2024. <a href="#">Synatra: Turning indirect knowledge into direct demonstrations for digital agents at scale</a> . In <i>Proceedings of NeurIPS</i> .	650
		651
		652
		653
		654
		655
	Jiayi Pan, Yichi Zhang, Nicholas Tomlin, Yifei Zhou, Sergey Levine, and Alane Suhr. 2024a. <a href="#">Autonomous evaluation and refinement of digital agents</a> . In <i>Proceedings of Conference on Language Modeling</i> .	656
		657
		658
		659
	Yichen Pan, Dehan Kong, Sida Zhou, Cheng Cui, Yifei Leng, Bing Jiang, Hangyu Liu, Yanyi Shang, Shuyan Zhou, Tongshuang Wu, and Zhengyang Wu. 2024b. <a href="#">Webcanvas: Benchmarking web agents in online environments</a> . In <i>Agentic Markets Workshop at ICML 2024</i> .	660
		661
		662
		663
		664
		665
	Victor Le Pochat, Tom van Goethem, Samaneh Tajalizadehkhoob, Maciej Korczynski, and Wouter Joosen. 2019. <a href="#">Tranco: A research-oriented top sites ranking hardened against manipulation</a> . In <i>Proceedings of Network and Distributed System Security Symposium</i> .	666
		667
		668
		669
		670
	Yujia Qin, Yining Ye, Junjie Fang, Haoming Wang, Shihao Liang, Shizuo Tian, Junda Zhang, Jiahao Li, Yunxin Li, Shijue Huang, et al. 2025. <a href="#">Ui-tars: Pioneering automated gui interaction with native agents</a> . <i>arXiv preprint arXiv:2501.12326</i> .	671
		672
		673
		674
		675
	Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy P Lillicrap. 2023. <a href="#">Androidinthewild: A large-scale dataset for android device control</a> . In <i>Proceedings of NeurIPS Datasets and Benchmarks Track</i> .	676
		677
		678
		679
		680
	Gaurav Sahu, Pau Rodríguez, Issam H. Laradji, Parmida Atighehchian, David Vázquez, and Dzmitry Bahdanau. 2022. <a href="#">Data augmentation for intent classification with off-the-shelf large language models</a> . In <i>Proceedings of the 4th Workshop on NLP for Conversational AI, ConvAI@ACL</i> .	681
		682
		683
		684
		685
		686
	Junhong Shen, Atishay Jain, Zedian Xiao, Ishan Amlekar, Mouad Hadji, Aaron Podolny, and Ameet Talwalkar. 2024. <a href="#">Scribeagent: Towards specialized web agents using production-scale workflow data</a> . <i>arXiv preprint arXiv:2411.15004</i> .	687
		688
		689
		690
		691
	Yu Su, Diyi Yang, Shunyu Yao, and Tao Yu. 2024. <a href="#">Language agents: Foundations, prospects, and risks</a> . In <i>Proceedings of EMNLP: Tutorial Abstracts</i> .	692
		693
		694
	Ruixiang Tang, Xiaotian Han, Xiaoqian Jiang, and Xia Hu. 2023. <a href="#">Does synthetic data generation of llms help clinical text mining?</a> <i>arXiv preprint arXiv:2303.04360</i> .	695
		696
		697
		698

699	Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024a. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. <i>arXiv preprint arXiv:2409.12191</i> .	Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022. <a href="#">Webshop: Towards scalable real-world web interaction with grounded language agents</a> . In <i>Proceedings of NeurIPS</i> .	757 758 759 760
707	Xingyao Wang, Yangyi Chen, Lifan Yuan, Yizhe Zhang, Yunzhu Li, Hao Peng, and Heng Ji. 2024b. <a href="#">Executable code actions elicit better LLM agents</a> . In <i>Proceedings of ICML</i> .	Jiacheng Ye, Jiahui Gao, Qintong Li, Hang Xu, Jiangtao Feng, Zhiyong Wu, Tao Yu, and Lingpeng Kong. 2022. <a href="#">Zerogen: Efficient zero-shot learning via dataset generation</a> . In <i>Proceedings of EMNLP</i> .	761 762 763 764
711	Zhiyong Wu, Chengcheng Han, Zichen Ding, Zhenmin Weng, Zhoumianze Liu, Shunyu Yao, Tao Yu, and Lingpeng Kong. 2024. <a href="#">OS-copilot: Towards generalist computer agents with self-improvement</a> . In <i>Proceedings of ICLR 2024 Workshop on Large Language Model (LLM) Agents</i> .	Ori Yoran, Samuel Joseph Amouyal, Chaitanya Malaviya, Ben Bogin, Ofir Press, and Jonathan Berant. 2024. <a href="#">Assistantbench: Can web agents solve realistic and time-consuming tasks?</a> In <i>Proceedings of EMNLP</i> .	765 766 767 768 769
717	Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, Yitao Liu, Yiheng Xu, Shuyan Zhou, Silvio Savarese, Caiming Xiong, Victor Zhong, and Tao Yu. 2024. <a href="#">OSWorld: Benchmarking multimodal agents for open-ended tasks in real computer environments</a> . In <i>Proceedings of NeurIPS Datasets and Benchmarks Track</i> .	Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao Liu, Yuxiao Dong, and Jie Tang. 2024. <a href="#">Agenttuning: Enabling generalized agent abilities for llms</a> . In <i>Findings of ACL</i> .	770 771 772 773
725	Nancy Xu, Sam Masling, Michael Du, Giovanni Campagna, Larry Heck, James A. Landay, and Monica Lam. 2021. <a href="#">Grounding open-domain instructions to automate web support tasks</a> . In <i>Proceedings of NAACL</i> .	Jiwen Zhang, Jihao Wu, Teng Yihua, Minghui Liao, Nuo Xu, Xiao Xiao, Zhongyu Wei, and Duyu Tang. 2024. <a href="#">Android in the zoo: Chain-of-action-thought for GUI agents</a> . In <i>Findings of EMNLP</i> .	774 775 776 777
730	Yiheng Xu, Dunjie Lu, Zhennan Shen, Junli Wang, Zekun Wang, Yuchen Mao, Caiming Xiong, and Tao Yu. 2024a. <a href="#">Agenttrek: Agent trajectory synthesis via guiding replay with web tutorials</a> . <i>arXiv preprint arXiv:2412.09605</i> .	Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. 2024. <a href="#">Gpt-4v(ision) is a generalist web agent, if grounded</a> . In <i>Proceedings of ICML</i> .	778 779 780
735	Yiheng Xu, Hongjin Su, Chen Xing, Boyu Mi, Qian Liu, Weijia Shi, Binyuan Hui, Fan Zhou, Yitao Liu, Tianbao Xie, Zhoujun Cheng, Siheng Zhao, Lingpeng Kong, Bailin Wang, Caiming Xiong, and Tao Yu. 2024b. <a href="#">Lemur: Harmonizing natural language and code for language agents</a> . In <i>Proceedings of ICLR</i> .	Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Yonatan Bisk, Daniel Fried, Uri Alon, et al. 2024. <a href="#">Webarena: A realistic web environment for building autonomous agents</a> . In <i>Proceedings of ICLR</i> .	781 782 783 784 785
742	Yiheng Xu, Zekun Wang, Junli Wang, Dunjie Lu, Tianbao Xie, Amrita Saha, Doyen Sahoo, Tao Yu, and Caiming Xiong. 2024c. <a href="#">Aguvis: Unified pure vision agents for autonomous gui interaction</a> . <i>arXiv preprint arXiv:2412.04454</i> .		
747	An Yan, Zhengyuan Yang, Wanrong Zhu, Kevin Lin, Linjie Li, Jianfeng Wang, Jianwei Yang, Yiwu Zhong, Julian McAuley, Jianfeng Gao, et al. 2023. <a href="#">Gpt-4v in wonderland: Large multimodal models for zero-shot smartphone gui navigation</a> . <i>arXiv preprint arXiv:2311.07562</i> .		
753	Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chunyuan Li, and Jianfeng Gao. 2023. <a href="#">Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v</a> . <i>arXiv preprint arXiv:2310.11441</i> .		

## Appendices

This supplementary material provides additional details omitted in the main text.

- Appendix A: Mind2Web Training and Evaluation Details
- Appendix B: Trajectory Synthesis Details
- Appendix C: More Related Work
- Appendix D: Reasoning Generation Agent
- Appendix E: Prompt Details
- Appendix F: Trajectory Examples

### A Mind2Web Training and Evaluation Details

Table A.2 shows the hyperparameters and training time for experiments on Mind2Web-Live and Multimodal-Mind2Web. All experiments use Nvidia H100 GPUs.

#### A.1 Mind2Web-Live

We exclude the following websites - <https://www.kbb.com>, <https://www.sixflags.com>, <https://www.viator.com>, <https://www.menards.com>, <https://www.amctheatres.com>, <https://www.cargurus.com>, <https://www.gamestop.com>, <https://www.cabelas.com>, <https://www.rei.com> due to denial of access faced during our tests. Table 5 shows the results on Mind2Web-Live for 83 out of 104 tasks across the remaining 37 websites. The results on the whole Mind2Web-Live evaluation set are given in Table A.1. The results in Table 5 are reported as the maximum over three runs, accounting for intermittent website access issues that may affect evaluation consistency. For Mind2Web-Live, the dataloader first samples training instances at the trajectory level and then randomly samples a step from the trajectory to construct the final training instance. Thus, the number of epochs is calculated at the trajectory level. We use a viewport resolution of  $1280 \times 720$  during inference. The Mind2Web-Live dataset is released under the MIT license, which permits its use in academic research.

#### A.2 Multimodal-Mind2Web

Following Deng et al. (2023), we obtain the top-50 elements from a pre-trained DeBERTa (He et al.,

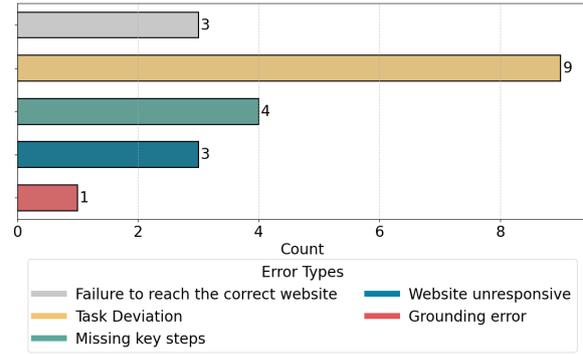


Figure A.1: Statistics for different error cases in Mind2Web-Live evaluation. Task deviation is the most prevalent error type.

2021) candidate generation model, which are then used to construct the accessibility tree and SoM image inputs. Following Ou et al. (2024), we always include the ground truth element in the input. We use a viewport resolution of  $1280 \times 720$  which includes the GT element during inference. We follow the setting in Zheng et al. (2024) and report element accuracy, operation F1, and step SR as evaluation metrics. All experiments on Multimodal-Mind2Web use a single training and evaluation run. The dataloader uniformly samples training instances from the set of action steps across all trajectories. The Multimodal-Mind2Web dataset is released under the Responsible AI license, which permits its use in academic research.

#### A.3 Ablation Studies

We conduct ablation studies to assess the impact of various design choices on overall performance (Table A.3). To evaluate the importance of visual modality, we experiment with using just the textual modality for the Phi-3.5V model, replacing it with the text-only Phi-3-mini (Abdin et al., 2024). In addition to Qwen2-VL-7B and Phi-3.5V, we also evaluate LLaVA-Mistral-7B (Liu et al., 2023), a strong MLLM baseline. Our results show that omitting the visual modality leads to a sharp 4.8% drop in performance for Phi-3.5V, underscoring its importance for effective GUI grounding. Furthermore, LLaVA-Mistral-7B significantly underperforms compared to both Qwen2-VL-7B and Phi-3.5V, highlighting the necessity of a stronger MLLM backbone for improved GUI agent performance.

#### A.4 Case Studies for Mind2Web-Live

We randomly sample 20 error cases for Explorer on Mind2Web-Live to gain insights for future im-

Model	Avg. Step SR (%)	Completion Rate (%)	Task SR (1) (%)	Full Task SR (%)
<b>API-based Models</b>				
GPT-4o	56.4	50.4	44.2	22.1
GPT-3.5	–	36.5	–	15.4
<b>Open-source Instructed Models</b>				
Mistral-7B-Instruct-0.3 (Jiang et al., 2023)	33.0	28.6	25.0	11.5
Qwen2-72B-Instruct (Bai et al., 2023)	–	<b>40.9</b>	–	15.4
Qwen2-VL-7B (Wang et al., 2024a)	37.9	33.3	31.7	12.5
Phi-3.5V (Abdin et al., 2024)	27.0	22.3	21.2	1.9
<b>Supervised Fine-Tuning</b>				
<b>Explorer-4B</b>	41.6	36.7	30.8	16.4
<b>Explorer-7B</b>	<b>42.0</b>	36.9	<b>32.7</b>	<b>16.4</b>

Table A.1: Results on Mind2Web-Live benchmark. The results for GPT-4, GPT-3.5, and Mistral-7B have been reproduced on our Linux servers. The full task success rate (SR) represents the successful completion of all key nodes for a given task. The average step success rate represents the proportion of completed key nodes, macro-averaged across tasks. The completion rate represents the proportion of completed key nodes, micro-averaged across tasks. Task SR (1) represents task SR with a tolerance of up to one error/key node. Our Phi-3.5V model, finetuned on synthetic trajectory data from Explorer, outperforms much larger models, including Mistral-7B and Qwen2-72B-Instruct, by a significant margin and is comparable to GPT-3.5.

Dataset	Model	Train Data	Hyperparameters	Train time (hours)
M2W-Live	Qwen2-VL-7B	Syn.	batch_size:64, epoch:2, learning_rate: $1 \times 10^{-5}$	15
	Qwen2-VL-7B	M2W	batch_size:64, epoch:2, learning_rate: $1 \times 10^{-5}$	1.5
	Qwen2-VL-7B	Syn. + M2W	batch_size:64, epoch:2, learning_rate: $1 \times 10^{-5}$	15.5
M2W-Live	Phi-3.5V	Syn.	batch_size:64, epoch:2, learning_rate: $4 \times 10^{-5}$	12.5
	Phi-3.5V	M2W	batch_size:64, epoch:2, learning_rate: $1 \times 10^{-5}$	1
	Phi-3.5V	Syn. + M2W	batch_size:64, epoch:2, learning_rate: $4 \times 10^{-5}$	12.5
Multi.-M2W	Qwen2-VL-7B	Syn.	batch_size:64, epoch:10, learning_rate: $4 \times 10^{-5}$	17
	Phi-3.5V	Syn.	batch_size:64, epoch:10, learning_rate: $4 \times 10^{-5}$	12

Table A.2: Hyperparameters used in our experiments.

Model	Avg. Step SR (%)	Completion Rate (%)	Full Task SR (%)
LLaVA-Mistral-7B	32.0	30.3	4.8
Phi-3-mini (text-only)	36.6	34.0	13.3
Phi-3.5V	44.0	39.4	18.1
Qwen2-VL-7B	<b>45.3</b>	<b>40.2</b>	<b>19.3</b>

Table A.3: Ablation studies on language models used for fine-tuning (Mind2Web-Live).

865 improvement. These errors fall into the following  
866 categories:

- 867 • *Task deviation*: The agent executes actions un-  
868 related to the given task, thus failing to com-  
869 plete it.
- 870 • *Missing key steps*: The agent retrieves results  
871 that partially satisfy the required constraints,  
872 e.g., the agent finds women’s clothes of the  
873 correct size but incorrect type or color.
- 874 • *Grounding error*: The agent fails to interact  
875 with a valid element on the page.
- 876 • *Website unresponsive*: The agent executes the  
877 correct action, but the website does not re-  
878 spond.
- 879 • *Failure to reach the correct website*: This hap-  
880 pens when the agent fails to output the correct  
881 website URL or use the search engine to arrive  
882 at the correct website.

883 Figure A.1 presents the statistics for these error  
884 types.

## 885 B Trajectory Synthesis Details

### 886 B.1 Cost Analysis

887 We use GPT-4o-turbo, which costs \$2.5 per 1M  
888 tokens for our trajectory synthesis. Each proposal  
889 or refinement stage uses 3.6K textual tokens on  
890 average. Each input image costs \$0.0028. The  
891 calculation assumes an average of 7.7 steps per  
892 trajectory, including the proposal stage. Table B.4  
893 shows the breakdown for the different stages of  
894 trajectory generation.

$$895 \text{ Total cost} = \$0.0128 * 7.7 + \$0.02581$$

$$896 \quad \quad \quad + \$0.02381 = \$0.148$$

897 The average cost per raw trajectory is \$0.15. The  
898 success rate is estimated as 53.1%. Thus, the aver-  
899 age cost per successful trajectory is estimated to be  
900 \$0.28.

Phase	Cost per step	Total cost
Proposal	\$0.0128	\$0.0128
Refinement	\$0.0128	\$0.0856
Verification	\$0.02381	\$0.02381
Summarization	\$0.02581	\$0.02581

Table B.4: Cost breakdown for different modules in the pipeline.

Action Type	Description	Count
click [elem]	Click on elem.	415K
type [elem] [text]	Type text	62K
select [elem] [text]	Select text from dropdown list.	5K
goto [url]	Go to url.	26K
search_google [query]	Search for query on Google.	4K
scroll [up/down]	Scroll up or down.	213K

Table B.5: Action space for web navigation in Explorer.

## B.2 Failure Modes of Trajectory Generation 901

902 We analyze cases where a generated trajectory is  
903 ultimately rejected by the task verifier agent. Our  
904 goal is to synthesize trajectory data that closely re-  
905 sembles human-annotated datasets for training web  
906 agents. However, since our pipeline collects tra-  
907 jectories through an exploration-driven approach,  
908 some trajectories result from random, incoherent  
909 action sequences that fail to align with a well-  
910 defined task intent. For instance, in shopping tasks,  
911 the agent may explore various products without  
912 demonstrating an intent to purchase (e.g., by adding  
913 items to the cart). Another failure case arises when  
914 the agent encounters errors on the final page due to  
915 automated browser detection, CAPTCHA verifica-  
916 tion, or an unresponsive website. We note that the  
917 verifier agent is instructed to judge a trajectory as  
918 successful if the task is completed, except for the  
919 final login and payment steps. The trajectories in  
920 failure modes are still valuable for web agents to  
921 learn low-level tasks such as form filling, basic in-  
922 teraction with web elements, and visual grounding.

## C More Related Work 923

### C.1 LLM-based Web Agents 924

925 Recent advances in multimodal language models  
926 have facilitated the development of web agents —  
927 autonomous systems designed to interact with real-  
928 world websites to perform everyday tasks (Deng  
929 et al., 2023; Hong et al., 2024; Cheng et al., 2024;  
930 Zheng et al., 2024). Web agents have made sig-  
931 nificant progress, evolving from simulated envi-  
932 ronments (Liu et al., 2018) to complex real-world  
933 applications (Deng et al., 2023; Yao et al., 2022;  
934 Zhou et al., 2024). Key challenges for web agents  
935 include long-term planning, visual grounding, and  
936 memory management. To improve long-context  
937 understanding, WebAgent (Gur et al., 2024) uti-  
938 lizes multiple LLMs - one for planning, summa-  
939 rization, and grounded program synthesis. SeeAct  
940 (Zheng et al., 2024) adopts a two-step procedure of  
941 planning followed by grounding at each step using  
942 GPT-4 to accomplish web agent tasks. Another line

943 of work employs a vision-only approach to train a  
944 GUI grounding model that directly predicts pixel  
945 coordinates for executing GUI agent tasks (Cheng  
946 et al., 2024; Kapoor et al., 2024; Gou et al., 2024).  
947 However, a significant bottleneck remains — the  
948 lack of large-scale, high-quality web trajectory data  
949 for training robust agents. Our work presents a new  
950 framework for synthesizing large-scale web trajec-  
951 tory data to train end-to-end web agents.

and Table E.11, respectively. The training prompt  
for Explorer is given in Table E.12.

990  
991

## 952 C.2 Web Agent Benchmarks and Datasets

953 Early benchmarks for web tasks such as Mini-  
954 Wob++ (Liu et al., 2018) focused on testing low-  
955 level actions on simulated websites. However,  
956 these simulated websites fail to capture the com-  
957 plexity of the real-world web. Mind2Web (Deng  
958 et al., 2023) introduces a trajectory-level dataset  
959 with 2K tasks across 137 real-world websites and  
960 31 domains. However, it employs a static evalu-  
961 ation method that penalizes alternative valid exe-  
962 cution paths. To overcome this limitation, follow-  
963 up work has explored alternative evaluation ap-  
964 proaches, including functional correctness-based  
965 evaluation in WebArena (Zhou et al., 2024) and  
966 key-node-based evaluation in Mind2Web-Live (Pan  
967 et al., 2024b). Towards the goal of making web  
968 agents more capable of performing realistic tasks,  
969 GAIA (Mialon et al., 2024) and AssistantBench  
970 (Yoran et al., 2024) introduce benchmarks that in-  
971 clude time-consuming information-seeking tasks.  
972 In this work, we develop Explorer, a multimodal  
973 web agent trained on our synthetic dataset, and  
974 showcase its strong performance across online  
975 and offline benchmarks, including Mind2Web-Live,  
976 Multimodal-Mind2Web, and MiniWob++.

## 977 D Reasoning Generation Agent.

978 Inspired by Xu et al. (2024c), the reasoning gener-  
979 ation agent is a pre-trained Qwen2-VL-7B model  
980 that takes as input the current action, high-level  
981 task description, screenshot, accessibility tree, and  
982 action history. It then outputs a post-hoc reasoning  
983 trace for performing that action. These reasoning  
984 traces are helpful for training GUI agents in a chain-  
985 of-thought style.

## 986 E Prompt Details

987 The prompts for the task proposer agent, task re-  
988 finer agent, task summarizer agent, and task verifier  
989 agent are given in Table E.7, Table E.9, Table E.10,

---

**System Role** What does this webpage show? Imagine you are a real user on this webpage. Given the webpage screenshot and parsed HTML/accessibility tree, please provide a single task that a user might perform on this page and the corresponding first action towards completing that task.

Do the following step by step:

1. Generate a single task that a user might perform on this webpage. Be creative and come up with diverse tasks
2. Given the webpage screenshot and parsed HTML/accessibility tree, generate the first action towards completing that task (in natural language form).
3. Given the webpage screenshot, parsed HTML/accessibility tree, and the natural language action, generate the grounded version of that action.

---

**ACTION SPACE:** Your action space is: ['click [element ID]', 'type [element ID] [content]', 'select [element ID] [content of option to select]', 'scroll [up]', 'scroll [down]', and 'stop'].

Action output should follow the syntax as given below:

'click [element ID]': This action clicks on an element with a specific ID on the webpage.

'type [element ID] [content]': Use this to type the content into the field with id. By default, the "Enter" key is pressed after typing. Both the content and the ID should be within square braces as per the syntax.

'select [element ID] [content of option to select]': Select an option from a dropdown menu. The content of the option to select should be within square braces. When you get (select an option) tags from the accessibility tree, you need to select the serial number (element\_id) corresponding to the select tag, not the option, and select the most likely content corresponding to the option as input.

'scroll [down]': Scroll the page down.

'scroll [up]': Scroll the page up.

---

**IMPORTANT:** To be successful, it is important to STRICTLY follow the below rules:

**Action generation rules:**

1. You should generate a single atomic action at each step.
2. The action should be an atomic action from the given vocabulary - click, type, select, scroll (up or down), or stop.
3. The arguments to each action should be within square braces. For example, "click [127]", "type [43] [content to type]", "scroll [up]", "scroll [down]".
4. The natural language form of action (corresponding to the field "action\_in\_natural\_language") should be consistent with the grounded version of the action (corresponding to the field "grounded\_action"). Do NOT add any additional information in the grounded action. For example, if a particular element ID is specified in the grounded action, a description of that element must be present in the natural language action.
5. If the type action is selected, the natural language form of action ("action\_in\_natural\_language") should always specify the actual text to be typed.
6. You should issue a "stop" action if the current webpage asks to log in or for credit card information.
7. To input text, there is NO need to click the textbox first, directly type content. After typing, the system automatically hits the 'ENTER' key.
8. STRICTLY Avoid repeating the same action (click/type) if the webpage remains unchanged. You may have selected the wrong web element.
9. Do NOT use quotation marks in the action generation.

---

**Task proposal rules:**

1. You should propose tasks that are relevant to the website and can be completed using the website.
2. You should only propose tasks that do not require login to execute the task.
3. You should propose tasks that are clear and specific.
4. For each task, provide concrete information or constraints, and use mock-up information (identifier, number, personal information, name, attributes, etc.) to make the task more specific and realistic.
5. The task description should provide all the necessary information to complete the task.
6. The task should be feasible to complete by a real user and should not require any additional information that is not available on the website.

The output should be in below format:

*Continued on next page*

*Continued from previous page*

---

	<b>OUTPUT FORMAT:</b> Please give a short analysis of the screenshot, parsed HTML/accessibility tree, then put your answer within ```, for example, "In summary, the proposed task and the corresponding action is: ```{"task": <TASK>:str, "action_in_natural_language":<ACTION_IN_NATURAL_LANGUAGE>:str, "grounded_action": <ACTION>:str}```
<b>User Role</b>	Website URL: {INIT_URL} Parsed HTML/Accessibility Tree: {A11Y_TREE} {SCREENSHOT}

---

Table E.7: Prompt for Task Proposer Agent.

---

**System Role** What does this webpage show? Imagine you are a real user on this webpage, and your overall task is {OVERALL\_TASK}. This is the list of actions you have performed that lead to the current page {PREV\_ACTION\_LIST}. You are also given the webpage screenshot and parsed HTML/accessibility tree.

Do the following step by step:

1. Please predict what action the user might perform next that is consistent with the previous action list in natural language.
2. Then based on the parsed HTML/accessibility tree of the webpage and the natural language action, generate the grounded action.
3. Update the overall task aligned with this set of actions.

---

**ACTION SPACE:** Your action space is: ['click [element ID]', 'type [element ID] [content]', 'select [element ID] [content of option to select]', 'scroll [up]', 'scroll [down]', and 'stop'].

Action output should follow the syntax as given below:

'click [element ID]': This action clicks on an element with a specific id on the webpage.

'type [element ID] [content]': Use this to type the content into the field with id. By default, the "Enter" key is pressed after typing. Both the content and the id should be within square braces as per the syntax.

'select [element ID] [content of option to select]': Select an option from a dropdown menu. The content of the option to select should be within square braces. When you get (select an option) tags from the accessibility tree, you need to select the serial number (element\_id) corresponding to the select tag, not the option, and select the most likely content corresponding to the option as input.

'scroll [down]': Scroll the page down.

'scroll [up]': Scroll the page up.

---

**IMPORTANT:** To be successful, it is important to STRICTLY follow the below rules:

**Action generation rules:**

1. You should generate a single atomic action at each step.
2. The action should be an atomic action from the given vocabulary - click, type, select, scroll (up or down), or stop
3. The arguments to each action should be within square braces. For example, "click [127]", "type [43] [content to type]", "scroll [up]", "scroll [down]".
4. The natural language form of action (corresponding to the field "action\_in\_natural\_language") should be consistent with the grounded version of the action (corresponding to the field "grounded\_action"). Do NOT add any additional information in the grounded action. For example, if a particular element ID is specified in the grounded action, a description of that element must be present in the natural language action.
5. If the type action is selected, the natural language form of action ("action\_in\_natural\_language") should always specify the actual text to be typed.
6. You should issue the "stop" action when the given list of input actions is sufficient for a web task.
7. You should issue a "stop" action if the current webpage asks to log in or for credit card information.
8. To input text, there is NO need to click the textbox first, directly type content. After typing, the system automatically hits the 'ENTER' key.
9. STRICTLY Avoid repeating the same action (click/type) if the webpage remains unchanged. You may have selected the wrong web element.
10. Do NOT use quotation marks in the action generation.

---

**Task proposal rules:**

1. You should propose tasks that are relevant to the website and can be completed using the website itself.
2. The overall task should be well-aligned to the entire set of actions in history plus the current generated action. It should not be focused just on the current action.
3. You should only propose tasks that do not require login to execute the task.
4. You should propose tasks that are clear and specific.
5. For each task, provide concrete information or constraints, and use mock-up information (identifier, number, personal information, name, attributes, etc.) to make the task more specific and realistic.
6. The task description should provide all the necessary information to complete the task.
7. The task should be feasible to complete by a real user and should not require any additional information that is not available on the website.

The output should be in below format:

*Continued on next page*

*Continued from previous page*

---

	<b>OUTPUT FORMAT:</b> Please give a short analysis of the screenshot, parsed HTML/accessibility tree, and history, then put your answer within ```, for example, "In summary, the proposed task and the corresponding action is: ```{"task": <TASK>:str, "action_in_natural_language": <ACTION_IN_NATURAL_LANGUAGE>:str, "grounded_action": <ACTION>:str}```
<b>User Role</b>	Website URL: {INIT_URL} Parsed HTML/Accessibility Tree: {A11Y_TREE} {SCREENSHOT}

---

Table E.9: Prompt for Task Refiner Agent.

---

<b>System Role</b>	<p>Given a list of actions performed on the website {WEBSITE_URL} and the corresponding screenshots  List of actions: {ACTION_LIST}  Your task is to come up with a single task description that will be accomplished by performing these actions in the given sequence on the website.</p> <hr/> <p><b>IMPORTANT:</b></p> <ol style="list-style-type: none"> <li>1. The task must contain some actions: “Buy, Book, Find, Check, Choose, show me, search, browse, get, compare, view, give me, add to cart, ...”, ideally involving transactions/finding information on a specific product or service.</li> <li>2. You should propose tasks that are clear and specific.</li> <li>3. The task description should provide all the necessary information to complete the task.</li> <li>4. The task description must indicate the domain of the website at the end of the task with the format: “... on task website”, for instance, “Purchase a laptop on Amazon”, “Book a hair appointment on Yelp”, etc.</li> <li>5. The task should be feasible to complete by a real user and should not require any additional information that is not specified in this input.</li> <li>6. The task description should specify constraints like given budget, product features, and other specifications that can narrow down the search to a particular item/product.</li> <li>7. Do NOT use any quotation marks (either single or double) in the task description.</li> </ol> <hr/> <p>The output should be in the below format:  <b>OUTPUT FORMAT:</b> Please first give some analysis of the actions and screenshots and then output the overall task description. put your answer within ```, for example, “In summary, the answer is: ```:&lt;TASK_DESCRIPTION&gt;:str```.</p>
--------------------	--

---

Table E.10: Prompt for Task Summarizer Agent.

---

<b>System Role</b>	<p>You are an expert in evaluating the performance of a web navigation agent. The agent is designed to help a human user navigate a website to complete a task. Given the user’s intent, the agent’s action history, the final state of the webpage, and the agent’s response to the user, your goal is to decide whether the agent’s execution is successful or not.  There are four types of tasks:</p> <ol style="list-style-type: none"> <li>1. <b>Transaction:</b> The user wants to perform a transaction on the webpage, such as booking a ticket, ordering a product, etc. The bot should at least initiate the add-to-cart or checkout process. It is still a success if the bot has done actions of ‘add to cart’ or checkout and encounters the login page. If the bot fails to do so, the task is considered a failure.</li> <li>2. <b>Information seeking:</b> The user wants to obtain certain information from the webpage, such as information of a product, reviews, map info, comparison of map routes, etc. The bot’s response must contain the information the user wants, or explicitly state that the information is not available. Otherwise, e.g. the bot encounters an exception and responds with the error content, the task is considered a failure. Besides, be careful about the sufficiency of the agent’s actions. For example, when asked to list the top-searched items in a shop, the agent should order the items by the number of searches, and then return the top items. If the ordering action is missing, the task is likely to fail.</li> <li>3. <b>Site navigation:</b> The user wants to navigate to a specific page. Carefully examine the bot’s action history and the final state of the webpage to determine whether the bot successfully completes the task. No need to consider the bot’s response.</li> <li>4. <b>Content modification:</b> The user wants to modify the content of a webpage or configuration. Carefully examine the bot’s action history and the final state of the webpage to determine whether the bot successfully completes the task. No need to consider the bot’s response.</li> </ol> <hr/> <p><b>IMPORTANT</b></p> <ul style="list-style-type: none"> <li>- If a product has been added to the bag/cart in the action list but just the purchase is pending, it should be counted as a success.</li> <li>- If you see the checkout page for the product you want to purchase, it should be counted as a success.</li> <li>- Format your response into two lines as shown below:</li> </ul> <p><b>Thoughts:</b> &lt;your thoughts and reasoning process&gt;  <b>Status:</b> "success" or "failure"</p>
--------------------	--

---

Table E.11: Prompt for Task Verifier Agent (adapted from Pan et al. (2024a)).

<b>System Role</b>	<p>You are an expert at completing instructions on webpage screens. You will be presented with a screenshot image with some numeric tags. If you decide to click somewhere, you should choose the numeric element index closest to the location you want to click. You should decide the action to continue this instruction. You will be given the accessibility tree of the current screen in the format: [element_idx] [role] [alt text or button name]. Here are the available actions:</p> <pre> {"action": "goto", "action_natural_language": str, "value": &lt;the URL to go to&gt;} {"action": "google_search", "action_natural_language": str, "value": &lt;search query for google&gt;} {"action": "click", "action_natural_language": str, "idx": &lt;element_idx&gt;} {"action": "type", "action_natural_language": str, "idx": &lt;element_idx&gt;, "value": &lt;the text to enter&gt;} {"action": "select", "action_natural_language": str, "idx": &lt;element_idx&gt;, "value": &lt;the option to select&gt;} {"action": "scroll [up]", "action_natural_language": str} {"action": "scroll [down]", "action_natural_language": str} </pre> <p>Your final answer must be in the above format.</p>
<b>User Role</b>	<p>The instruction is to {TASK DESCRIPTION}.  History actions: {PREVIOUS ACTIONS}  Here is the screen information: {ACCESSIBILITY TREE}  Think about what you need to do with the current screen, and output the action in the required format in the end.</p>

Table E.12: Prompt for Web Agent Training.

## **F Trajectory Examples**

Figure F.2 shows a sample trajectory executed on the IKEA website. Figure F.3 shows the set-of-mark annotations and accessibility tree inputs of the model during trajectory generation, training, and inference.

Task description: Navigate to the IKEA US website and browse to find three-seat fabric sofas

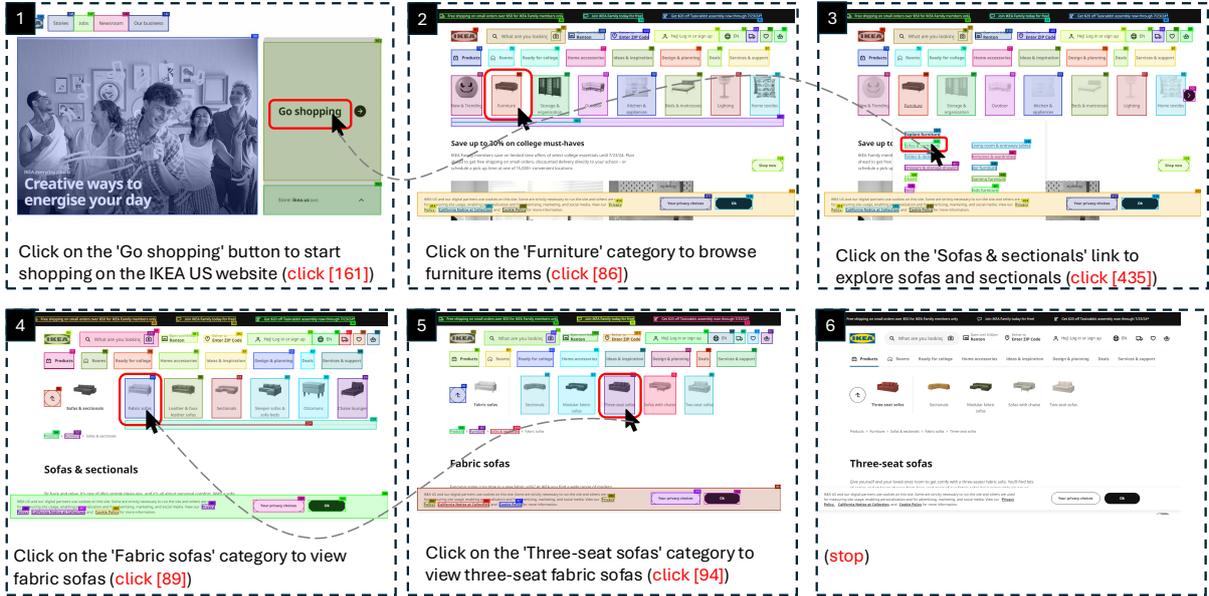
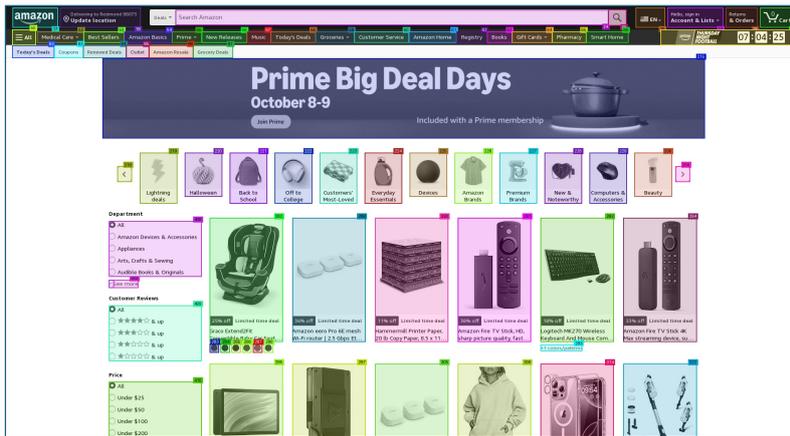


Figure F.2: Example synthetic trajectory from Explorer. Each step shows the set-of-mark annotated screenshot along with the grounded action taken by the GPT-4 agent.



(a) Set-of-mark annotated screenshot of webpage

- [17] [A] [Amazon]
- [18] [A] [Delivering to Redmond 98073 Update location]
- [24] [INPUT, TYPE=TEXT] [Search Amazon]
- [25] [INPUT, TYPE=SUBMIT] []
- [26] [A] [Choose a language for shopping.]
- [27] [A] [Hello, sign in Account & Lists]
- [28] [A] [Returns & Orders]
- [29] [A] [0 items in cart]
- [50] [A] [Open Menu]
- [52] [A] [Medical Care]
- [53] [A] [Best Sellers]
- [54] [A] [Amazon Basics]
- [55] [A] [Prime]
- [56] [A] [New Releases]
- ...
- [289] [A] [36% off Limited time deal Amazon eero Pro 6E mesh Wi-Fi router | 2.5 Gbps Ethernet | Coverage up to 6,000 sq. ft | Connect 100+ devices | Ideal for streaming, working, and gaming | 3-Pack | 2022 release Amazon eero Pro 6E mesh Wi-Fi router | 2.5 Gbps Etu2026]
- ...
- [408] [DIV] [Department]
- [409] [DIV] [Customer Reviews]
- [410] [DIV] [Price]\*\*,

(b) Corresponding A11y tree

Figure F.3: Visualization of the model inputs during trajectory generation, model training, and inference. The example corresponds to step 2 of the trajectory in Figure 1.