

# AGENTREK: AGENT TRAJECTORY SYNTHESIS VIA GUIDING REPLAY WITH WEB TUTORIALS

Anonymous authors

Paper under double-blind review

## ABSTRACT

Graphical User Interface (GUI) agents hold great potential for automating complex tasks across diverse digital environments, from web applications to desktop software. However, the development of such agents is hindered by the lack of high-quality, multi-step trajectory data required for effective training. Existing approaches rely on expensive and labor-intensive human annotation, making them unsustainable at scale. To address this challenge, we propose *AgentTrek*, a scalable data synthesis pipeline that generates high-quality web agent trajectories by leveraging web tutorials. Our method automatically gathers tutorial-like texts from the internet, transforms them into task goals with step-by-step instructions, and employs a visual-language model (VLM) agent to simulate their execution in a real digital environment. A VLM-based evaluator ensures the correctness of the generated trajectories. We demonstrate that training GUI agents with these synthesized trajectories significantly improves their grounding and planning performance over the current models. Moreover, our approach is more cost-efficient compared to traditional human annotation methods. This work underscores the potential of guided replay with web tutorials as a viable strategy for large-scale GUI agent training, paving the way for more capable and autonomous digital agents.

## 1 INTRODUCTION

Graphical User Interfaces (GUIs) are a fundamental medium for human-computer interaction, enabling users to perform tasks across various digital platforms. Automating GUI operations through agentic automation has the potential to significantly enhance productivity by enabling autonomous task completion using human-centric tools. Additionally, this approach can foster the development of advanced AI systems capable of learning from rich digital environments.

Recent advancements in large language models (LLMs) have endowed the models with powerful abilities in understanding, reasoning, and decision-making, which are essential for the evolution of GUI agents in diverse contexts such as web (Zheng et al., 2024), desktop (Xie et al., 2024), and mobile applications (Zhang et al., 2023). Despite these advancements, the performance of GUI agents remains suboptimal. Contemporary Large Language Models (LLMs) are primarily engineered and trained on datasets optimized for generating informative responses (Ouyang et al., 2022; OpenAI, 2024). Their architecture and training paradigms are not inherently designed to make complex, sequential action decisions that require long-term observation and historical context. Consequently, training GUI agents with multi-step trajectory data is crucial to improving their capabilities.

High-quality agent trajectories contain several key components: a high-level goal, a sequence of interleaved observations, natural language reasoning, and grounded actions (as shown in Figure 1).

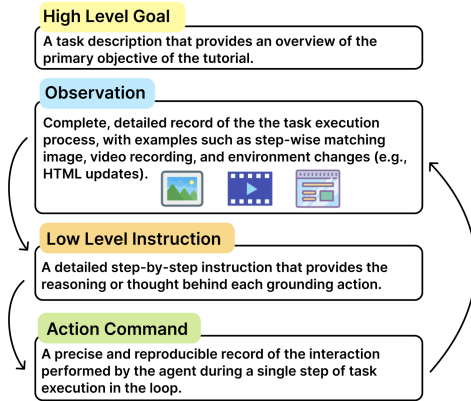


Figure 1: Expected Agent Trajectories

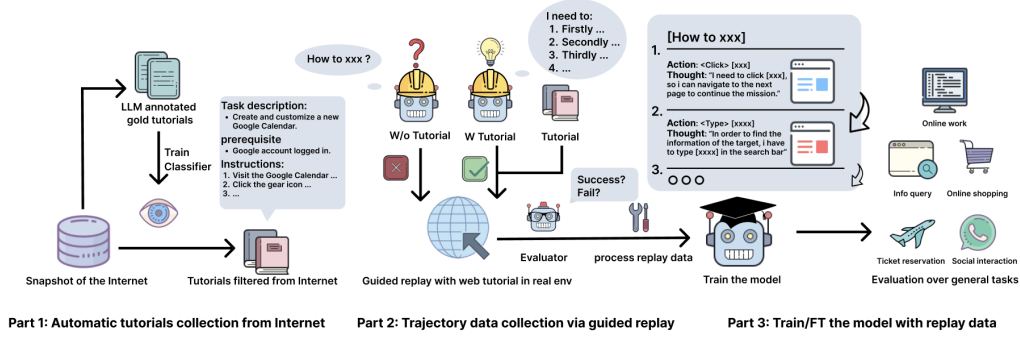


Figure 2: **Overview of the AgentTrek Pipeline:** (1) *Automatic Tutorial Collection from the Internet:* Tutorial-related data is extracted and filtered from internet sources using heuristic methods and a FastText model. An LLM processes the filtered textual data, transforming it into structured tutorials. (2) *Trajectory data collection via guided replay:* A VLM agent interacts with the real digital environment guided by tutorials, while high-quality trajectory data, including observations, actions, and reasoning, is collected. Another VLM evaluator acts as a judge to further improve the effectiveness of the synthetic dataset. (3) *Training and fine-tuning with replay data:* The collected trajectory data is used to train and fine-tune GUI agent models, which are evaluated on standard agent benchmarks, demonstrating significant improvements.

Unfortunately, such data is not readily available on the internet like textual or image data, as it involves complex situational reasoning and multimodal interactivity. Existing approaches typically rely on human annotation to collect these trajectories (Deng et al., 2024; Rawles et al., 2023; Li et al., 2024), a process that is both expensive and not scalable.

To address this data scarcity, data synthesis has emerged as a vital approach in AI system development. Synthesizing agent trajectories presents significant challenges due to the need for interwoven natural language instructions, visual observations, and context-specific actions that must be accurately grounded in the GUI environment. Although there have been some successful applications of LLMs in data synthesis pipelines (Ye et al., 2022; Peng et al., 2023; Qin et al., 2023), these complexities still make GUI trajectory synthesis particularly demanding.

In this work, we present AgentTrek, a scalable data synthesis pipeline specifically designed for training GUI agents. We begin by automatically gathering and filtering tutorial-like text from the web, which describes GUI tasks and workflows in web environments. These tutorials are then transformed into agent tasks with high-level objectives and detailed step-by-step instructions. Using a visual-language model (VLM) agent, we simulate the execution of these tasks, guided by the synthesized tutorials. An evaluator model is also employed to subsequently verify whether the goal was successfully achieved. Through this comprehensive pipeline, we efficiently generated a large volume of high-quality web agent trajectories.

Our experimental results demonstrate that training GUI agent models with these synthesized web agent trajectories not only improves their performance but also enables them to surpass the capabilities of their initial teacher models, which is the replay model GPT-4 in our case. Compared to traditional human-annotated data pipelines, our method is significantly more cost-effective, emphasizing the scalability and economic viability of the AgentTrek pipeline.

- We introduce AgentTrek, a novel pipeline that leverages web tutorials to synthesize high-quality web agent trajectory data at scale, effectively bridging the gap between LLM capabilities and the demanding need for multi-step, context-rich training data for GUI agents.
- Extensive experiments demonstrate that agents trained with our synthesized data outperform those trained on existing datasets in both grounding and planning capabilities, validating the effectiveness of AgentTrek.
- Our pipeline significantly reduces the cost and scalability obstacles of human-annotated data collection, providing a practical approach for large-scale GUI agent training through data synthesis.

Table 1: **Comparison of AgentTrek with other trajectory datasets for training.** For the calculation of dataset size and average steps, see Appendix A.

Datasets	Size	Average Steps	HTML	AxTree	Intermediate Reasoning	Video	Matching Screenshot	Website	Task Inst. Level
RUSS	80	5.4	Yes	No	No	No	No	22	Low
ScreenAgent	203	4.3	No	No	Yes	No	Yes	-	High & Low
WebLINX	969	18.8	Yes	No	No	No	Yes	155	High & Low
MM-Mind2Web	1009	7.3	Yes	No	No	No	No	137	High
GUIAct	2482	6.7	No	No	No	No	Yes	121	High
<b>AgentTrek (Ours)</b>	<b>4902</b>	<b>12.1</b>	<b>Yes</b>	<b>Yes</b>	<b>Yes</b>	<b>Yes</b>	<b>Yes</b>	<b>127</b>	<b>High &amp; Low</b>

## 2 METHOD

We introduce a pipeline to collect and process GUI tutorials from the internet for training visual language models (VLMs) in web automation tasks. The method comprises three main steps:

1. **Collecting Tutorials:** We extract web interaction tutorials from large datasets using keyword filtering and language models to identify and standardize relevant content.
2. **Guided Replay:** An agent uses these tutorials to perform tasks in a web environment, interacting with real websites while we record its actions and thoughts.
3. **Model Training:** We train a visual agent model that relies on screenshots and standard GUI actions, enhancing its web navigation capabilities with the collected data.

This approach enables efficient training of VLMs without extensive manual annotation, offering a scalable solution for automating web tasks.

### 2.1 AUTOMATIC TUTORIALS COLLECTION FROM INTERNET

We first extract web interaction tutorials from Redpajama dataset (Computer, 2023). A rule-based heuristic filter is applied to create a preliminary dataset, a subset of which is annotated by an advanced LLM to generate labeled samples for training a effective FastText classification model (Joulin et al., 2017), the tutorial classifier. This classifier further enhances the data quality through filtering. Finally, LLMs are employed to tag and paraphrase the raw text of tutorials into a standardized format, preparing them for the replay phase in Section 2.2.

#### 2.1.1 PREFILTER FUNCTION

Although GUI tutorials are abundant online, they constitute only a small fraction of web content, making a pre-filter essential for identifying relevant content. Similar patterns often appear in tutorials, such as distinctive keywords like ‘click’ and ‘type’, as well as platform-specific terms like ‘macOS’ and ‘Windows’. We compiled a rule-based filter using keyword lists sourced from official websites and forums. Leveraging RedPajama data with over 20 billion URLs, our pre-filter applies **Keyword Matching** in the first 38k words, evaluates samples based on **Length**, and filters them by **URL Format** for relevance.

Validated using 180 positive and 105 negative ground-truth samples, the prefilter achieved a 92.69% recall rate on positive samples, ensuring both diversity and quantity. After filtering, the dataset size is reduced from 20.8 billion to 68.8 million entries (Figure 4).

#### 2.1.2 LLM LABELER

While initial rule-based filtering narrows the context, the proportion of true positive tutorial content remains low. To improve the quality and relevance of selected tutorials, we leverage an advanced LLM, GPT-4o MINI, for automated labeling, due to its strong ability to comprehend and analyze complex, information-dense text. Prior to full implementation, we tested GPT-4o MINI on a manually annotated ground-truth validation set, where it achieved an F1 score nearly 90%. In cases where human and LLM annotations conflicted, the LLM demonstrated the ability to identify tutorial-related

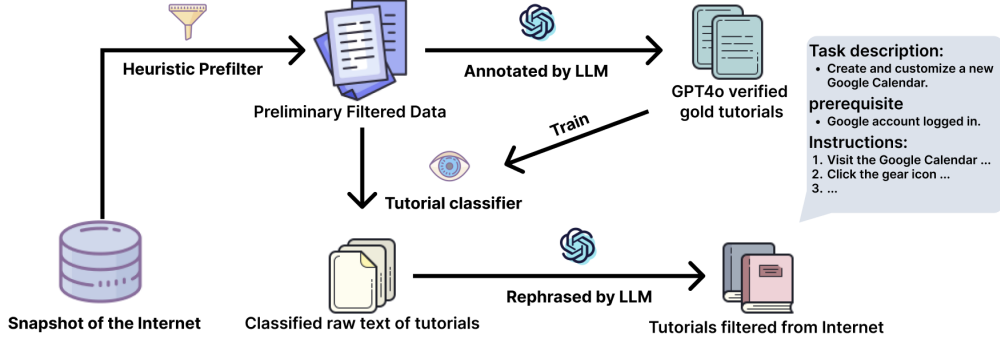


Figure 3: **Overview of the tutorial filtering and classification pipeline.** Starting with Redpajama, the data is prefiltered, annotated by an advanced LLM, and used to train a tutorial classifier. The classifier further filters the raw text, which is then paraphrased into structured tutorials with task descriptions, prerequisites, and step-by-step instructions.

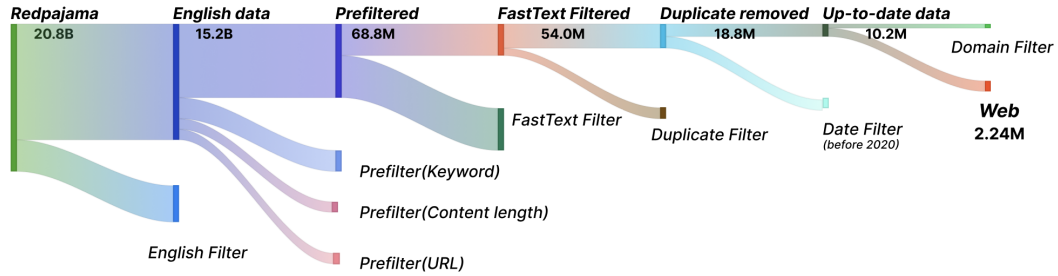


Figure 4: The data flow during the early stages of our pipeline.

content in lengthy texts that humans might overlook. This, along with the validation set result, suggests that GPT-4O MINI may surpass human performance in webpage labeling, enabling efficient generation of a large labeled dataset for training in the following section.

### 2.1.3 FASTTEXT FILTER

Following the automated labeling process, we employed FastText, an n-gram-based deep learning model, as our classifier. FastText classifies tutorial text segments as tutorial or non-tutorial, with a binary output and a confidence score to enhance the accuracy of tutorial selection. To train the model, we combined LLM-labeled data with human-labeled samples, creating a dataset of approximately 90,000 examples. The train-test split is 95:5, with the model demonstrating strong classification performance. Using this classifier, we further curated the initial filtered dataset, collecting approximately 18.8 million tutorial-like web text samples.

Table 2: Performance of Filters.

Metric	Precision	Recall	F1
<b>Prefilter</b>	0.69	0.61	0.60
<b>LLM</b>	0.885	0.885	0.89
<b>FastText</b>	0.895	0.895	0.89

### 2.1.4 TAG & PARAPHRASE

After filtering the raw tutorial content using the FastText model, we then tag and paraphrase the content for further processing, including extracting meta-information and formatting the tutorials according to a standardized template. To handle the length and noise of the raw tutorial data, we em-

ployed GPT-4O MINI, streamlining the tagging and paraphrasing while ensuring the output aligned with the comprehensive template and gold-standard examples.

The key components of the template include specifying the **Platform and Target** (e.g., macOS, Windows, browser, or app), providing a concise **Task Description**, listing **Prerequisites** needed before starting, outlining **Step-by-Step Instructions** for completing the task, and detailing the **Expected Outcome**. The cost for tagging and paraphrasing 1,000 entries is approximately 0.89 dollars.

## 2.2 TRAJECTORY DATA COLLECTION VIA GUIDED REPLAY

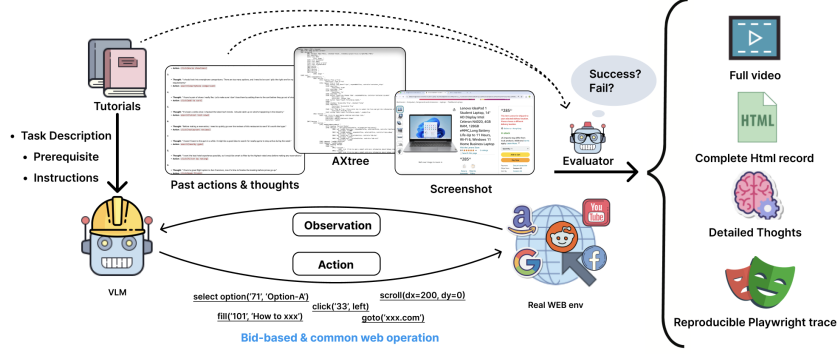


Figure 5: **Overview of Guided Replay data collection and evaluation pipeline.** A VLM agent is provided with the filtered and formatted tutorials, then observes and interacts with the real environment during execution, while all the actions and intermediate thoughts are recorded as data trajectory. The final result is evaluated by an advanced VLM to ensure the correctness.

### 2.2.1 TRAJECTORY DATA DEFINITION

The trajectory data generated by our pipeline is designed to enhance an agent’s web navigation capabilities by integrating high-level planning, low-level instructions, and grounded operations. Each data instance includes the following components:

**Task Information.** Detailed task metadata, including platform, task description, prerequisites, instructions, and expected outcomes, which support both planning and execution.

**Post-processed Textual Trajectory.** Refined after replay, highlighting key elements for model fine-tuning. This includes *Task Metadata*, summarizing the task to encourage adaptive decision-making, *Observations* to provide visual context, *Intermediate Reasoning* offering insights into the agent’s decision-making process, and *Action Sequence* to capture detailed element information for web interactions.

**Screenshots and Video Recordings.** Visual records of the entire process for comprehensive documentation.

**Reproducible Native Trace.** Captured via Playwright, including DOM snapshots, HTML, network flow, and action sequences, allowing full reconstruction and detailed analysis of agent-environment interactions.

### 2.2.2 GUIDED REPLAY WITH TUTORIALS

Although we have collected and processed high-quality tutorials, a significant gap remains in acquiring the grounding data crucial for training a more effective agent model. To address this, we leverage BrowserGym (Drouin et al., 2024) to enable the model to replay tasks under the guidance of the generated tutorials.

BrowserGym is a versatile environment for web task automation in the Chromium browser, enabling Visual Language Model (VLM) agents to execute web-based operations (Drouin et al., 2024). Agents are provided with tagged and paraphrased tutorials and a *target\_web\_url*, allowing them to



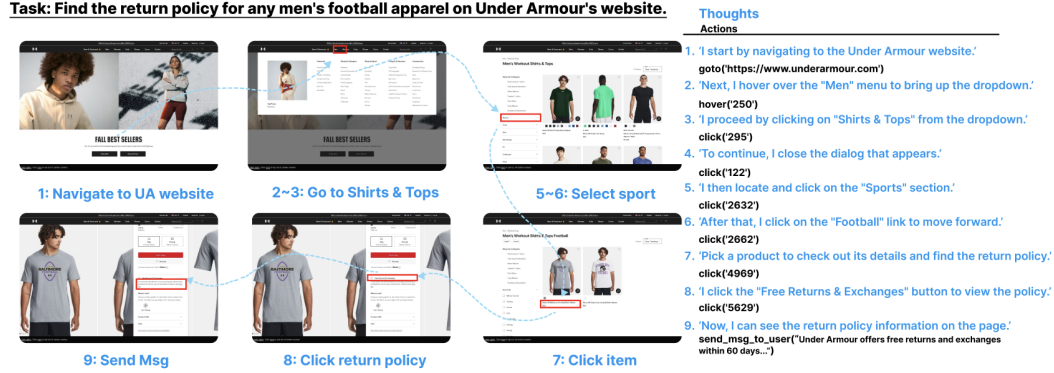


Figure 6: **Guided replay example.** This example demonstrates an agent’s execution of finding the return policy for men’s football apparel, showcasing its actions alongside the corresponding inner thoughts.

navigate directly to the task scene. Step-by-step instructions guide the agent through the task, with the expected outcome determining task completion.

The agent’s initial observations include the webpage’s viewport screenshot and accessibility tree (AXTree), but the HTML file is excluded due to its size and irrelevance to visual agents. Actions are executed using [Playwright](#) (Microsoft, 2023) functions such as `click`, `select_option`, and `clear`, while Playwright also records detailed traces, including target elements, coordinates, screenshots, and DOM snapshots at the same time, along with agent’s internal thoughts between actions.

Token consumption is about 8,027 per step and 86,114 per task. With GPT-4o-08-06, replaying 1,000 tasks costs approximately 215 dollars. Cost detail see C

### 2.2.3 EVALUATION OF TRAJECTORY

Although a large amount of guided replay data has been recorded, it is crucial to extract the effective segments that can truly contribute to enhancing the agent’s performance. Recent work by (Pan et al., 2024) highlights the potential of Vision Language Models (VLMs) in evaluating trajectory data using recorded images and interaction processes as input. VLMs are highly scalable, capable of processing large datasets concurrently at a low cost, and provide transparent evaluations. Therefore, we implemented a VLM Evaluator to further improve our data quality.

**VLM Evaluator Design.** To ensure trajectory data quality, we define *effectiveness* based on two criteria: adherence to task instructions and successful completion of core components. We employ GPT-4o as the backbone of our VLM evaluator, using a structured prompt to assess recorded trajectories. The evaluator receives the task description  $d$ , the agent’s action history  $\bar{a}$ , and inner thoughts  $\bar{r}$  for each step. The sequential format is: {task description; inner thought 1; action 1; inner thought 2; action 2; ...}, as illustrated in Figure 5. The VLM provides a trajectory-level assessment and performs stepwise analysis, offering justifications for any ineffective trajectory and identifying the earliest point of failure.

Table 3: Accuracy Comparison

Trajectory	Evaluator	Acc.
Web Tutorials	VLM Eval	84.0%
	GPT-4V	80.6%
WebArena	Cap. + GPT-4	82.1%
	Cap. + Mixtral	74.4%

Table 4: Cost Breakdown

Phase	Cost/1k (\$)	Model
T&P	0.89	mini
Replay	215.36	08-06
Eval	3.10	08-06
<b>Total</b>	<b>219.35</b>	—

**Validation on Human-Annotated Set.** Although the capabilities of Vision Language Models (VLMs) are well-recognized, validation is essential. To assess the automatic evaluator’s perfor-

mance, we manually reviewed 1,081 trajectories and created a validation set of 558 samples with human-annotated justifications.

As shown in the Table 3, despite handling different input formats across various task scenarios, the evaluator achieved high performance metrics. And according to the observation shown in Appendix D, evaluator often applies stricter standards than human evaluators. This demonstrates its robustness in accurately identifying effective trajectories.

## 2.3 TRAIN AND FINE-TUNE THE MODEL WITH TRAJECTORY DATA

We chose a purely visual agent model that relies exclusively on screenshot-based observations, rather than incorporating accessibility trees or textual representations, for several reasons. First, GUIs are inherently visual, and mapping instructions to visual elements aligns more closely with human cognitive processes. Second, Textual representations, such as HTML or accessibility trees, are often verbose, which leads to heavy overhead for computation. Second, Different websites can have varying structures for their textual representation while image-based representations allow the model to unify its observations across diverse platforms with varying resolutions, improving generalization.

### 2.3.1 PURE VISION & GUI ACTION FRAMEWORK

In this work, we propose to unify observation and action space via pure vision and standard pyautogui commands with a pluggable action system.

Using pure vision as input eliminates the need for the model to understand different UI source codes across platforms, even when visually similar elements are written differently in HTML. Additionally, HTML input typically costs an average of 4,000 tokens per step. In contrast, recent VLMs with high-resolution multimodal understanding, such as Qwen2-VL, require only 1,200 tokens for a 720p image. This significantly lowers the computational cost while maintaining sufficient visual information for the task.

For action, we choose the widely used standard pyautogui action space with a pluggable action system. Most web agent leverage playwright action space. But playwright actions incline to interact with html selector element instead of visual ui element. Therefore, we use pyautogui commands to unify basic GUI operations on web. Since we collect the data from website by playwright, we need to map the playwright actions to pyautogui actions as shown in the Figure 8. In addition, we utilize a pluggable action system to cover specific playwright action like select\_option.

### 2.3.2 MODEL ARCHITECTURE AND TRAINING

Unlike agents that rely on structured UI representations like accessibility trees, vision-based grounding requires models to map intents directly to visual observations. For this, we chose Qwen2-VL (Wang et al., 2024), which uses NaViT as an image encoder with dynamic resolution support (Dehghani et al., 2023). By removing absolute position embeddings and incorporating 2D-RoPE (Su et al., 2024), Qwen2-VL can process images of any resolution, efficiently converting them into variable visual tokens. This makes Qwen2-VL ideal for GUI agents, as it can encode high-resolution images with fewer token costs, making it well-suited for our tasks.

Our training process, starting with a VLM capable of high-resolution image understanding, consists of one tuning stage. We use data from AgentTrek Data to enhance VLM capabilities in grounding and planning.

## 3 EXPERIMENTS

AgentTrek automatically collects thousands of trajectories with detailed information, including inner thoughts and precise action coordinates, offering an ideal foundation for fine-tuning VLM into a reliable visual GUI agent. A successful visual GUI agent requires two key abilities: planning and grounding. We conducted extensive experiments to validate that AgentTrek data enhances the agent’s capabilities in both areas.

### 3.1 EXPERIMENTAL SETUP

To validate the effectiveness of our dataset, it is essential to demonstrate its impact on improving both the grounding and planning capabilities of the model. Therefore, we will evaluate the model’s performance on benchmarks that assess these abilities.

For grounding ability, we employ the ScreenSpot benchmark (Cheng et al., 2024), which evaluates the model’s single-step GUI grounding performance across various platforms. To assess planning ability, we utilize the Multimodal-Mind2Web benchmark (Deng et al., 2024), the multimodal extension of the web agent benchmark Mind2Web (Deng et al., 2024).

### 3.2 WEB GROUNDING

ScreenSpot is a GUI visual grounding benchmark comprising 1.2K single-step instructions and target element bounding boxes. It spans mobile, desktop, and web environments, categorizing elements into text and icons. Given that our data originates exclusively from the web, we focus solely on web-based performance.

Fine-tuning with the AgentTrek dataset significantly improved Qwen2-VL’s grounding ability for both text and icon-based tasks, more than doubling its baseline performance and surpassing several models on the ScreenSpot benchmark. This demonstrates the strong impact of AgentTrek in enhancing the model’s grounding capabilities for web-based GUI tasks.

Table 5: Comparison of grounding performance on ScreenSpot Web Grounding

Model	Text	Icon/Widget	Average
GPT-4 (Cheng et al., 2024)	9.2	8.8	9.0
GPT-4o (Cheng et al., 2024)	12.2	7.8	10.1
Qwen2-VL	35.2	25.7	30.7
SeeClick (Cheng et al., 2024)	55.7	32.5	44.7
CogAgent (Cheng et al., 2024)	70.4	28.6	50.7
GPT-4 + OmniParser (Lu et al., 2024)	81.3	51.0	67.0
Qwen2-VL w/ <b>AgentTrek</b>	<b>81.7</b>	<b>51.5</b>	<b>67.4</b>

### 3.3 WEB PLANNING

The baseline Qwen2-VL model is excluded due to its poor performance in locating target elements, a key requirement for web-based tasks. Thus, only the fine-tuned versions are presented in the table.

Fine-tuning with the AgentTrek dataset significantly improved Qwen2-VL’s performance, particularly in the Operation F1 metric, where it outperformed both GPT-3.5 and GPT-4 across all settings.

The combination of AgentTrek and Mind2Web datasets yields the best performance across all metrics and settings. Although the model achieve great performance after fine-tuning with mind2web dataset, it can still benefits from the AgentTrek Data.

These results highlight the complementary nature of the two datasets: AgentTrek provides high-quality grounded data with precise coordinates, while Mind2Web offers valuable data for handling complex web-based tasks. [More explanation about result source in Appendix J.2.](#)

## 4 ANALYSIS

With our AgentTrek pipeline, we generate large-scale trajectory data that excels in three areas. First, the dataset offers extensive diversity, covering multiple domains and task types, and benefiting from internet-sourced tutorials that enhance task execution. Our experiment showed a 230% performance increase when agents followed detailed instructions. Second, the data is gathered from real-world web environments, avoiding the limitations of simulations. Starting with RedPajama, we filtered and processed 12,526 tutorials, producing 4,902 successful trajectories from 127 websites.



Table 6: Performance comparison across different methods and evaluation settings. 'H', 'I', 'AT', 'M2W' stand for HTML, Image, AgentTrek, Mind2Web

Obs	Model	Method	Cross-Task			Cross-Website			Cross-Domain		
			Ele.Acc	Op.F1	Step SR	Ele.Acc	Op.F1	Step SR	Ele.Acc	Op.F1	Step SR
HTML	GPT-3.5	Choice	19.4	59.2	16.8	14.9	56.5	14.1	25.2	57.9	24.1
	GPT-4	Choice	40.8	63.1	32.3	30.2	61.0	27.0	35.4	61.9	29.7
H + I	GPT-4	Choice	46.4	73.4	40.2	38.0	67.8	32.4	42.4	69.3	36.8
	GPT-4	SoM	29.6	-	20.3	20.1	-	13.9	27.0	-	23.7
Image	Qwen2-VL										
	+ AT	Vision	45.5	84.9	40.9	40.8	82.8	35.1	48.6	84.1	42.1
	+ M2W	Vision	54.8	89.5	50.9	52.9	83.9	44.9	51.8	86.8	47.7
	+ AT + M2W	Vision	60.8	88.9	55.7	57.6	88.1	51.4	56.0	87.5	52.6

Third, the data is comprehensive, capturing high- and low-level task details, including DOM/HTML structures, AXTree snapshots, video recordings, and screenshots. This rich data improves the agent’s performance on long-horizon tasks, and with a per-trajectory cost of just \$0.551, our pipeline offers an efficient, scalable solution for data generation.

#### 4.1 IMPORTANCE OF TUTORIALS

Tutorials extracted from the internet play a crucial role in guiding the replay process. First, they ensure diversity in the generated trajectories. Tutorials often have distinct task goals, and even when they target the same objective, they may offer different execution methods, enriching the trajectory data. Second, tutorials significantly improve the agent’s execution. We tested the agent on 400 tasks, replaying them twice: once with tutorials and once using only high-level goals. The results show that step-by-step instructions greatly enhanced performance. Without tutorials, only 63 effective trajectories were generated (15.78% of the total). With tutorials, the agent produced 208 effective trajectories (52%), marking an increase of over 230%, demonstrating the importance of detailed instructions in improving reliability and effectiveness. For further analysis, please see [Appendix B](#)

#### 4.2 DATA COMPOSITION

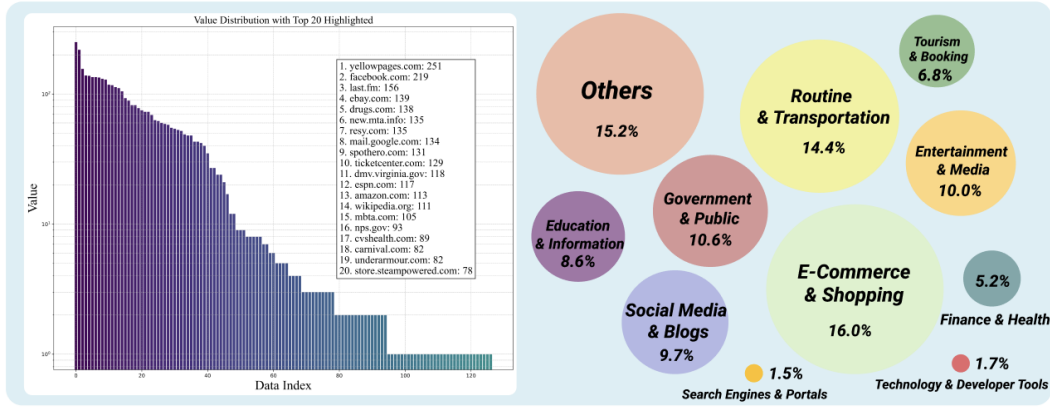


Figure 7: The distribution of website with domains involved in our dataset

To summarize the data flow through our pipeline: First, we filter tutorial data from the RedPajama web snapshot. Next, the filtered data is paraphrased for clarity and classification. We then gather up-to-date data from mainstream websites for replay and, finally, collect effective trajectory data from the replays.

After filtering RedPajama’s vast dataset, we retained over 18.8 million entries. By applying criteria such as recency and popularity, 12,526 tutorials were prepared for replay. With a success rate of 39.1%, we generated 4,902 trajectories, covering 127 websites across 11 distinct categories.

### 4.3 COMPARISON WITH EXISTING WORK AND RESEARCH CHALLENGES

AgentTrek generates comprehensive, large-scale trajectory data, excelling in several key areas as shown in Table 1 (Niu et al., 2024; Lù et al., 2024; Deng et al., 2024; Yao et al., 2022; Song et al., 2024; Wornow et al., 2024). First, with nearly 5k verified trajectories and an average of 12.1 steps per trajectory, our dataset provides a strong foundation for training and evaluating agents on long-horizon web tasks. Second, it is the most complete dataset to date, incorporating DOM/HTML structures, AXTree data, intermediate reasoning steps, full video recordings, and corresponding screenshots for each action.

Moreover, despite full automation without human intervention, our dataset maintains diversity across 120 websites and 12 distinct task categories. By leveraging modern large language models (LLMs), we can extract both high-level task objectives and detailed step-by-step instructions, offering flexibility for future use.

Finally, our pipeline significantly reduces the cost and scalability challenges of human-annotated data collection. With a success rate factored in, the cost per trajectory is just \$0.551, making our approach both efficient and scalable for large-scale data generation. Cost detail see C

## 5 RELATED WORK

**LLM-based Agents.** LLM-based agents are autonomous systems that leverage large language models (Brown et al.) to interact with real-world websites and os environments. These agents can understand natural language instructions and perform a wide range of complex tasks across various domains, such as e-commerce, online assistance, and knowledge navigation (Nakano et al., 2021; Cheng et al., 2024). Recent efforts in this space include models like SeeAct (Zheng et al., 2024) and WebVoyager (He et al., 2024), which aim to generalize agent behavior to real-world websites. While LLM-based agents have shown promise, challenges remain in the need for agent specified data. Our work extends this line of research by introducing a cost-effective pipeline to generate comprehensive agent trajectory data, advancing the state-of-the-art in data synthesis for agent-based applications.

**Agent Data.** As agents gain increasing popularity, the demand for efficient and scalable data is becoming both larger and more urgent. However, most existing data primarily serve as supplements to various benchmarks (Zhou et al., 2023; Li et al., 2023; Deng et al., 2024), with few datasets specifically designed for agent trajectory analysis. Furthermore, these datasets are often limited by the need for human annotation, which hampers scalability. In our work, our pipeline managed to automatically generate comprehensive agent trajectory data in a cost-efficient manner, paving the way for a new direction in data synthesis within the field of agents.

**Automatic Evaluation for Digital Agents.** Recently, there has been growing interest in automating the evaluation of digital agents using Vision-Language Models (VLMs) and Large Language Models (LLMs). These methods leverage models to assess agent performance in real-world tasks. Research in this area spans several dimensions: some works focus on trajectory-level success (Pan et al., 2024), while others evaluate stepwise success based on adherence to instructions (Wornow et al., 2024). Additionally, evaluations are conducted across various task environments, such as web-based platforms and mobile operating systems like Android and iOS (Pan et al., 2024). In our work, we prompt a VLM, GPT-4o, as an autonomous evaluator, using agent’s interaction process as inputs to assess whether the agent has successfully completed tasks at the trajectory level.

## 6 CONCLUSION

In this work, we introduce AgentTrek, an efficient pipeline designed to automatically generate comprehensive and cost-effective agent trajectory data. Additionally, we present a large and diverse dataset generated using this approach, which we validate by training models and evaluating their performance with promising result. Our research establishes a novel and promising direction for the future development of LLM agent, particularly in the automatic and low-cost synthesis of trajectory data. AgentTrek serves as a strong standard for agent data generation, setting the stage for future advancements in this field.

## REFERENCES

- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners.
- Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Yantao Li, Jianbing Zhang, and Zhiyong Wu. Seeclick: Harnessing gui grounding for advanced visual gui agents. *arXiv preprint arXiv:2401.10935*, 2024.
- Together Computer. Redpajama: an open dataset for training large language models, 2023. URL <https://github.com/togethercomputer/RedPajama-Data>.
- Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su. Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing Systems*, 36, 2024.
- Alexandre Drouin, Maxime Gasse, Massimo Caccia, Issam Hadj Laradji, Manuel Del Verme, Tom Marty, L’eo Boisvert, Megh Thakkar, Quentin Cappart, David Vazquez, Nicolas Chapados, and Alexandre Lacoste. Workarena: How capable are web agents at solving common knowledge work tasks? *ArXiv*, abs/2403.07718, 2024. URL <https://api.semanticscholar.org/CorpusID:268363855>.
- Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Yong Dai, Hongming Zhang, Zhenzhong Lan, and Dong Yu. Webvoyager: Building an end-to-end web agent with large multimodal models. *arXiv preprint arXiv:2401.13919*, 2024.
- Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. Bag of tricks for efficient text classification. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pp. 427–431. Association for Computational Linguistics, April 2017.
- Hanyu Lai, Xiao Liu, Iat Long Iong, Shuntian Yao, Yuxuan Chen, Pengbo Shen, Hao Yu, Hanchen Zhang, Xiaohan Zhang, Yuxiao Dong, and Jie Tang. Autowebglm: A large language model-based web navigating agent, 2024. URL <https://arxiv.org/abs/2404.03648>.
- Tao Li, Gang Li, Zhiwei Deng, Bryan Wang, and Yang Li. A zero-shot language agent for computer control with structured reflection. *arXiv preprint arXiv:2310.08740*, 2023.
- Wei Li, William Bishop, Alice Li, Chris Rawles, Folawiyo Campbell-Ajala, Divya Tyamagundlu, and Oriana Riva. On the effects of data scale on computer control agents, 2024. URL <https://arxiv.org/abs/2406.03679>.
- Yadong Lu, Jianwei Yang, Yelong Shen, and Ahmed Awadallah. Omniparser for pure vision based gui agent, 2024. URL <https://arxiv.org/abs/2408.00203>.
- Xing Han Lù, Zdeněk Kasner, and Siva Reddy. Weblinx: Real-world website navigation with multi-turn dialogue, 2024.
- Microsoft. Playwright for python documentation. <https://playwright.dev/python/>, 2023.
- Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- Runliang Niu, Jindong Li, Shiqi Wang, Yali Fu, Xiyu Hu, Xueyuan Leng, He Kong, Yi Chang, and Qi Wang. Screenagent: A vision language model-driven computer control agent. 2024.
- OpenAI. Gpt-4v(ision) system card, 2024. URL <https://openai.com/research/gpt-4v-system-card>.

- Tianyue Ou, Frank F. Xu, Aman Madaan, Jiarui Liu, Robert Lo, Abishek Sridhar, Sudipta Sengupta, Dan Roth, Graham Neubig, and Shuyan Zhou. Synatra: Turning indirect knowledge into direct demonstrations for digital agents at scale, 2024. URL <https://arxiv.org/abs/2409.15637>.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke E. Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Francis Christiano, Jan Leike, and Ryan J. Lowe. Training language models to follow instructions with human feedback. *ArXiv*, abs/2203.02155, 2022. URL <https://api.semanticscholar.org/CorpusID:246426909>.
- Jiayi Pan, Yichi Zhang, Nicholas Tomlin, Yifei Zhou, Sergey Levine, and Alane Suhr. Autonomous evaluation and refinement of digital agents. *ArXiv*, abs/2404.06474, 2024. URL <https://api.semanticscholar.org/CorpusID:269009430>.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with gpt-4. *ArXiv*, abs/2304.03277, 2023. URL <https://api.semanticscholar.org/CorpusID:257985497>.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou, Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. Toolllm: Facilitating large language models to master 16000+ real-world apis, 2023. URL <https://arxiv.org/abs/2307.16789>.
- Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy Lillicrap. Android in the wild: A large-scale dataset for android device control, 2023. URL <https://arxiv.org/abs/2307.10088>.
- Yifan Song, Da Yin, Xiang Yue, Jie Huang, Sujian Li, and Bill Yuchen Lin. Trial and error: Exploration-based trajectory optimization for llm agents. 2024.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Ke-Yang 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. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *ArXiv*, abs/2409.12191, 2024. URL <https://api.semanticscholar.org/CorpusID:272704132>.
- Michael Wornow, Avani Narayan, Ben Viggiano, Ishan S. Khare, Tathagat Verma, Tibor Thompson, Miguel Angel Fuentes Hernandez, Sudharsan Sundar, Chloe Trujillo, Krish Chawla, Rongfei Lu, Justin Shen, Divya Nagaraj, Joshua Martinez, Vardhan Agrawal, Althea Hudson, Nigam H. Shah, and Christopher Re. Do multimodal foundation models understand enterprise workflows? a benchmark for business process management tasks. *arXiv preprint arXiv:2406.13264*, 2024. URL <https://hazyresearch.stanford.edu/wonderbread-website>.
- Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, et al. Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments. *arXiv preprint arXiv:2404.07972*, 2024.
- Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chun yue Li, and Jianfeng Gao. Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v. *ArXiv*, abs/2310.11441, 2023. URL <https://api.semanticscholar.org/CorpusID:266149987>.
- Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents. *Advances in Neural Information Processing Systems*, 35:20744–20757, 2022.
- Jiacheng Ye, Jiahui Gao, Quintong Li, Hang Xu, Jiangtao Feng, Zhiyong Wu, Tao Yu, and Lingpeng Kong. Zerogen: Efficient zero-shot learning via dataset generation. *ArXiv*, abs/2202.07922, 2022. URL <https://api.semanticscholar.org/CorpusID:246867045>.

Chi Zhang, Zhao Yang, Jiaxuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. Appagent: Multimodal agents as smartphone users, 2023.

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.

Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, et al. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023.



## Table of Contents in Appendix

<b>A</b>	<b>Calculation of Other Trajectory Datasets</b>	<b>15</b>
<b>B</b>	<b>Analysis of the Effectiveness of Tutorials</b>	<b>15</b>
<b>C</b>	<b>Cost Details</b>	<b>15</b>
<b>D</b>	<b>Evaluator Alignment</b>	<b>15</b>
<b>E</b>	<b>Action Mapping</b>	<b>16</b>
<b>F</b>	<b>Experimental Results on Textual Data</b>	<b>16</b>
<b>G</b>	<b>Details in Collecting Tutorials</b>	<b>16</b>
	G.1 Prefilter Function . . . . .	16
	G.2 LLM Labeler Prompt . . . . .	18
	G.3 Tagging & Paraphrasing Prompt and Format . . . . .	18
<b>H</b>	<b>Examples of Failed Guided Replay Trajectories</b>	<b>20</b>
<b>I</b>	<b>Scaling up AgentTrek</b>	<b>21</b>
<b>J</b>	<b>Evaluation Benchmarks</b>	<b>21</b>
	J.1 GUI Grounding Evaluation . . . . .	21
	J.2 Offline GUI Agent Evaluation . . . . .	21

## A CALCULATION OF OTHER TRAJECTORY DATASETS

- **RUSS:** Cited based on the data provided in the table from WebLINX (Lù et al., 2024).
- **ScreenAgent:** Statistics obtained from the dataset available at <https://github.com/niuzaisheng/ScreenAgent/tree/main/data/ScreenAgent/train>.
- **WebLINX:** Calculated based on the train set information from Table 8 in (Lù et al., 2024) and data on HuggingFace (excluding the "say" actions), resulting in a total of 18,249 non-say actions with 969 demos.
- **Mind2Web:** Statistics derived from <https://huggingface.co/datasets/osunlp/Mind2Web>, specifically from the training subset.
- **Webshop (agent-eto):** Data statistics sourced from <https://huggingface.co/datasets/agent-eto/eto-sft-trajectory>.
- **WonderBread:** Calculations based on data presented in (Wornow et al., 2024).

## B ANALYSIS OF THE EFFECTIVENESS OF TUTORIALS

Key factors contributing to this improvement include:

1. **Direct Access to Target URL:** Tutorials provide the target URL, allowing direct access to the initial task state, reducing errors in locating the correct webpage.
2. **Assisted Planning with Human Expertise:** Tutorials aid in planning by providing steps informed by human experience, which tend to be reliable, thereby reducing the likelihood of errors during task execution and bridging the gap in the agent’s knowledge for unknown tasks.
3. **Navigating Multi-Level Menus:** Tutorials offer clear paths to hidden elements, preventing the agent from failing due to incorrect navigation through complex menus.

## C COST DETAILS

In this part we provide the details of our cost in generating trajectory data with via our pipeline:

Phase	Cost per 1,000 Entries (USD)	Model Used
Tag and Paraphrase	0.886	gpt-4o-mini
Replay	215.359	gpt-4o-2024-08-06
Evaluator	3.104	gpt-4o-2024-08-06

Table 7: Cost breakdown for each phase in the process

Another two important factors are the ratio of web-related tutorials (0.275) and the Replay Success Rate (39.9%). Using these, we can calculate the cost per verified effective trajectory as follows:

$$\text{Cost per trajectory} = \frac{\text{Tag and Paraphrase price}}{\text{Web ratio}} + \frac{\text{Replay price} + \text{Evaluate price}}{\text{Replay Success Rate}}$$

The cost per 1,000 verified effective trajectories is 550.75 \$.

## D EVALUATOR ALIGNMENT

In this part, we provide the details of metrics between the human and automatic evaluator.

Trajectory	Evaluator	Accuracy
Web Tutorials	VLM Evaluator	84.0%
Webarena	GPT-4V	80.6%
	Captioner + GPT-4	82.1%
	Captioner + Mixtral	74.4%

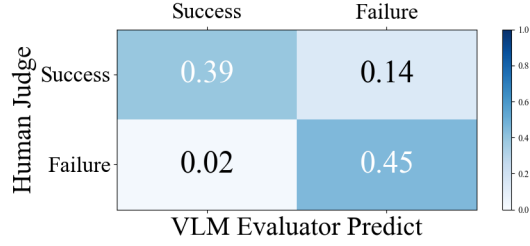


Figure 8: Confusion Matrix of our VLM evaluator’s performance on the human-annotated validation set, compared with evaluators across different scenarios.

## E ACTION MAPPING

Table 8: Mapping between Playwright and PyAutoGUI Action Spaces.

Category	Playwright Action	PyAutoGUI Action
Basic Actions	<code>page.click()</code>	<code>pyautogui.click()</code>
	<code>page.type()</code>	<code>pyautogui.write()</code>
	<code>page.press()</code>	<code>pyautogui.press()</code>
	<code>page.hover()</code>	<code>pyautogui.moveTo()</code>
	<code>page.scroll()</code>	<code>pyautogui.scroll()</code>
Advanced Actions	<code>page.fill()</code>	<code>pyautogui.write()</code> (clearing)
	<code>page.dblclick()</code>	<code>pyautogui.doubleClick()</code>
	<code>page.dragAndDrop()</code>	<code>pyautogui.dragTo()</code>
	<code>page.clear()</code>	<code>pyautogui.click()</code> <code>pyautogui.hotkey(ctr1, A)</code> <code>pyautogui.press(delete)</code>
Plugin	<code>playwright.select_option()</code>	<code>browser.select()</code>

## F EXPERIMENTAL RESULTS ON TEXTUAL DATA

To provide further supports for the effective of our AgentTrek data, we conducted an experiment to evaluate the performance of a pure textual agent using the textual version data of our AgentTrek trajectories. This allow us to study the contribution of textual modalities of AgentTrek.

We fine-tuned the Qwen2.5-7B-Instruct model using AgentTrek trajectories that included accessibility tree as observations and playwright actions as the agent’s action space. We then evaluated the model on WebArena, an OOD web agent benchmark featuring self-hosted websites. These websites are entirely out-of-domain (OOD) from the AgentTrek dataset, ensuring that the evaluation reflects the model’s generalization capability.

We fine-tuned **Qwen2.5-7B-Instruct** on AgentTrek’s textual data and achieved the following results on WebArena and Miniwob++ as shown in Table 9. We observe that our fine-tuned model achieves the highest performance among open-source web agents and approaches the performance of GPT-4o, demonstrating the effectiveness of AgentTrek data in improving real-world web agent capabilities and generalization across modalities.

## G DETAILS IN COLLECTING TUTORIALS

### G.1 PREFILTER FUNCTION

- **Keyword Density:** The web content must contain a minimum of 20 common keywords, ensuring sufficient topic coverage.
- **Keyword Diversity:** The text must incorporate at least 4 distinct common keywords.

**System Prompt**

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 task goal, the agent's trajectory, your goal is to decide whether the agent's execution is successful or not.

**\*Evaluation Criteria\***

Whether the agent's trajectory is effective and corresponding to the goal

**\*Instructions\***

1. Review the agent's actions and reasoning processes step by step.
2. if the agent is stuck in the very first login stage, which means it fails to log into target website at the beginning, that's a failure.
3. Determine if the agent has achieved the task goal based on the trajectory. A task can be considered successful if most trajectory is effective.
4. the agent sometimes can't stop after finishing a task and continue doing repeated actions. these actions may be some failed attempt after a series of correct actions. the task should be regarded as successful if the correct actions are effective and almost reach the goal.
5. if the agent is stuck in the loop at the early stage of the task, which means they don't even get close to the goal before they get stuck in the loop, that's a failure. for example, the agent begin to get stuck before third step.
6. when the task is to change the google account password, it can't be regarded as successful when agent finish at trying to click "manage your account".
7. if there are over 8 correct action in the trajectory, it can be regard as a successful agent.
8. final saving action is not a must. the task is successful if the agent does most things right and just forget to save the change at last.
9. if the original task has 2 subtasks, the agent only complete one of them, that's still a success. e.g. the task is to update name and birthday, but agent only update name, that's fine.
10. if the task is to post a review, the agent can be considered successful when it finish writing the review and reach the step to post it, don't have to click the post button.
11. Since we don't have a printer, some printing related task can be considered successful if the agent reach the step to click print button.
12. if the task is finished at the initial state and the agent do nothing because of it, it should also be regarded as successful.

**\*IMPORTANT\***

1. in the trajectory, an action always follows a corresponding reasoning, which shows the observation and thought of the agent.
2. your response should be contain:  
Thoughts: <your thoughts and reasoning process>  
Status: "success" or "failure"

**User Prompt**

The goal of the task: {task.des}  
trajectory: {trajectory}

Figure 9: Prompts to query the VLM Autonomous Evaluator.

- **Essential Keyword Frequency:** At least one mandatory keywords must appear multiple times (minimum twice) within the content, demonstrating topic relevance.

Table 9: Comparison of task success rate on WebArena

Model	WebArena	Miniwob++
CodeLlama-7B-Instruct (Ou et al., 2024)	0.00	23.04
LLaMa3-chat-8B (Ou et al., 2024)	3.32	31.74
Qwen2.5-7B-Instruct	3.80	30.19
LLaMa3-chat-70B (Ou et al., 2024)	7.02	48.70
GPT-4o(Zhou et al., 2023)	13.10	-
GPT-4(Ou et al., 2024)	14.41	53.04
Synatra-CodeLlama-7B (Ou et al., 2024)	6.28	38.20
AutoWebGLM (OOD SFT) (Lai et al., 2024)	8.50	-
AutoWebGLM (In-domain RFT) (Lai et al., 2024)	18.20	-
Qwen2.5-7B-Instruct w/ <b>AgentTrek</b>	<b>10.46</b>	45.28

## G.2 LLM LABELER PROMPT

To achieve more precise and context-aware labeling, we designed the following prompt to guide the LLM in further assessing whether the given URL and context meet our requirements, as illustrated in Fig 10.

<b>System Prompt</b>
You are an assistant that classifies content based on specific criteria. Your task is to evaluate whether a given piece of content serves as a tutorial specifically related to graphical user interfaces (GUI), such as for web applications, desktop applications, or operating systems.
<b>Classification Criteria</b>
The content qualifies as a GUI-related tutorial if it meets the following conditions:
1. It includes a task description outlining what needs to be achieved.
2. It provides clear step-by-step instructions for interacting with a GUI, such as:
- Step 1: Open the application
- Step 2: Navigate to the settings menu
Given the URL and context, determine if the content is a GUI-related tutorial or not. Output '1' if it is a GUI-related tutorial and '0' if it is not. Provide only the number as the output.
<b>User Prompt</b>
- URL: {url}
- Context: {context}

Figure 10: User Prompt for Classifying GUI Tutorials

## G.3 TAGGING & PARAPHRASING PROMPT AND FORMAT

Here we present the prompt designed to utilize LLM to help do the tagging & paraphrasing of the identified GUI-tutorial related context.



**User Prompt**

The following is a tutorial from the website. It may contain several tutorials. Please extract the first tutorial only and format the first tutorial according to the specified schema:

**Text:** {context}

**Schema:**


```
{
  "platform":
    "Platform category (choose from: macOS, Windows (Default if not
specified in the tutorial), Linux, Android, iOS)",
  "target type":
    "Type of platform (choose from: Web browser, PC app, Mobile app,
PC operating system, Mobile operating system, where the tutorial's
steps are performed). Tutorials that involve interacting with the
browser software itself, such as 'opening Chrome settings,' should be
classified as a PC app type.",
  "target object":
    "Specific name of the web browser or (non web browser)
applications or operating system where the tutorial's steps are
performed (e.g., Chrome browser (Default for browser and web
tutorial), Microsoft Excel (app name), Windows system settings)",
  "target web URL":
    "The exact URL of the web page where the tutorial's actions take
place, applicable only if the target object is a web browser (e.g.,
None, https://mail.google.com, https://www.amazon.com,
https://github.com). Be careful, the URL provided at the beginning
is always not the URL where the tutorial's actions are about. For
example, a tutorial from
https://abidakon.com/how-to-make-google-slide-vertical/ about
changing Google Slides, its target web URL should be
https://docs.google.com/presentation.",
  "task description":
    "Task description text (Provide a concise summary in one
sentence, including essential details)",
  "prerequisites":
    "Prerequisite text describing necessary conditions before
starting the task",
  "instructions":
    [
      "Step_1: Instruction text describing the action to be taken",
      // Following instructions
    ]
  "instructions steps":
    "Total number of instructions steps",
  "expected result":
    "Text describing the expected result after following the
instructions"
}
```

Figure 11: User Prompt for Extracting and Formatting GUI Tutorials

## H EXAMPLES OF FAILED GUIDED REPLAY TRAJECTORIES

**Task** Find and reserve a table at a top restaurant in Boston for April 22.

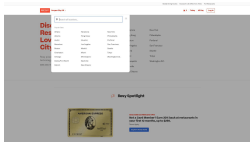
**Step 1**



**Thought:** I have successfully navigated to <https://resy.com>. Now, I need to follow the instructions to find and reserve a table at a top restaurant in Boston for April 22. The first step is to click on the location field and type 'Boston'.

**Action:** click('52')

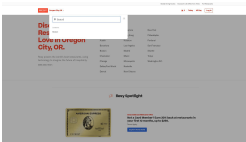
**Step 2**



**Thought:** I have successfully clicked on the location field button (bid 52). Now, I need to enter 'Boston' into this field as per Step 2 of the instructions. The location field should now be active and ready for input.

**Action:** fill('61', 'Boston')


**Step 3**



**Thought:** I have successfully entered 'Boston' into the location field. The next step is to select the date for the reservation, which should be April 22. I need to find and click on the date field to open the calendar picker.

**Action:** click('159')

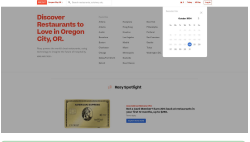
**Step 4**



**Thought:** I have clicked on the date field button (bid '159') to open the calendar for selecting the reservation date. The calendar should now be visible and expanded. I need to locate the button for April 22 and click on it to set the desired date for the reservation.

**Action:** click('523')


**Step 5**



**Thought:** I have opened the date selection calendar, but it's currently showing October 2024. I need to navigate back to April 2024 to select April 22 for the reservation. The calendar doesn't seem to have a direct way to jump to a specific month, so I'll need to use the 'Previous month' button multiple times to reach April 2024.

**Action:** click('517')

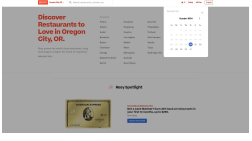
**Step 6**



**Thought:** I previously clicked the 'Next month' button to navigate through the calendar. However, I need to go back 18 months from October 2024 to reach April 2024. Clicking the 'Next month' button was not the correct action. I should look for a 'Previous month' button instead.

**Action:** click('606')


**Step 7**



**Thought:** The previous action to click the 'Previous month' button with bid '606' resulted in an error, indicating that the element could not be found. This suggests that the button's bid may have changed or the calendar view has been updated. I need to reassess the current state of the page and find an alternative way to navigate to April 2024.

**Action:** click('516')

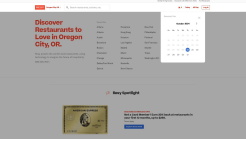
**Step 8**



**Thought:** Based on the error message, it seems that the 'Previous month' button with bid '516' is not available on the current page. This could be because the calendar view has changed or the button's bid has been updated. I need to reassess the available elements and find an alternative way to navigate to April 2024.

**Action:** click('510')

**Step 9**



**\*\*FAIL\*\***

Figure 12: Replay agent was unable to complete booking before the actual date due to the tutorial expiration.

Table 10: Multimodal Mind2Web Step SR across varying amounts of training data from AgnetTrek (AT).

Data Amount	Cross-Task	Cross-Website	Cross-Domain
20%	36.1	35.5	39.5
40%	41.0	35.8	42.5
60%	41.6	37.2	42.8
80%	42.6	38.0	44.3
100%	42.6	37.5	45.0

## I SCALING UP AGENTTREK

In this section, we further scale up the data amount of AgnetTrek (more than 10K trajectories) to explore the effectiveness of AgentTrek in large-scale size. We trained the model using varying proportions of the dataset (20% to 100%) and assessed its performance on Multimodal-Mind2Web across three splits. The results are presented in Table 10. We observe that the performance improves steadily as more data is used, with the best results achieved when using the full dataset. This underscores the value of scaling up AgentTrek in improving model effectiveness.

## J EVALUATION BENCHMARKS

In this section, we introduce more details of evaluation benchmarks used in our work.

### J.1 GUI GROUNDING EVALUATION

**ScreenSpot.** ScreenSpot (Cheng et al., 2024) is a benchmark developed specifically for GUI visual grounding tasks, featuring 1.2K single-step instructions along with the coordinates of target elements. The dataset includes diverse grounding instructions tailored for mobile, desktop, and web platforms and categorizes elements into text and icons/widgets. Two distinct assessment scenarios are utilized: (1) *Original Instructions*, where models directly execute grounding actions as per the provided instructions; and (2) *Self-plan*, where models are expected to formulate plans in natural language based on the original instructions before carrying out the grounding actions.

### J.2 OFFLINE GUI AGENT EVALUATION

**Multimodal-Mind2Web.** We evaluated the offline planning capabilities of GUI agents on websites using the Multimodal-Mind2Web benchmark (Zheng et al., 2024), which is an extension of the original Mind2Web benchmark (Deng et al., 2024). Performance was measured using Element Accuracy (Ele.Acc), Operation F1 (Op.F1), and Step Success Rate (Step SR).

The GPT-3.5 and GPT-4 results for the HTML and HTML+Image observation are derived from the SeeAct (Zheng et al., 2024) method. For the Choice method, it employs a DeBERTa-base cross-encoder to rank the interactive elements on the current HTML page. The top 50 elements are selected as options, and the GPT-3.5/GPT-4 model then chooses one of these elements as the answer. For the SoM method (Yang et al., 2023; Zheng et al., 2024), it renders a new webpage image by adding red bounding boxes and labels to every HTML node in the source code, allowing GPT-4 to understand the webpage screenshot and identify the target action object by referring to the labels. For the Image observation, as detailed in Section 2.3.1, we use Qwen2VL (Wang et al., 2024) within a pure vision framework. Specifically, Qwen2VL processes only the webpage screenshots and specifies the target action object by generating its coordinates.