

# WebVoyager : Building an End-to-End Web Agent with Large Multimodal Models

Hongliang He<sup>1,3\*</sup>, Wenlin Yao<sup>2</sup>, Kaixin Ma<sup>2</sup>, Wenhao Yu<sup>2</sup>, Yong Dai<sup>2</sup>,  
Hongming Zhang<sup>2</sup>, Zhenzhong Lan<sup>3</sup>, Dong Yu<sup>2</sup>

<sup>1</sup>Zhejiang University, <sup>2</sup>Tencent AI Lab, <sup>3</sup>Westlake University

hehongliang@westlake.edu.cn, wenlinyao@global.tencent.com

## Abstract

The rapid advancement of large language models (LLMs) has led to a new era marked by the development of autonomous applications in real-world scenarios, which drives innovation in creating advanced web agents. Existing web agents typically only handle one input modality and are evaluated only in simplified web simulators or static web snapshots, greatly limiting their applicability in real-world scenarios. To bridge this gap, we introduce WebVoyager, an innovative Large Multimodal Model (LMM) powered web agent that can complete user instructions end-to-end by interacting with real-world websites. Moreover, we establish a new benchmark by compiling real-world tasks from 15 popular websites and introduce an automatic evaluation protocol leveraging multimodal understanding abilities of GPT-4V to evaluate open-ended web agents. We show that WebVoyager achieves a 59.1% task success rate on our benchmark, significantly surpassing the performance of both GPT-4 (All Tools) and the WebVoyager (text-only) setups, underscoring the exceptional capability of WebVoyager. The proposed automatic evaluation metric achieves 85.3% agreement with human judgment, indicating its effectiveness in providing reliable and accurate assessments of web agents.<sup>1</sup>

## 1 Introduction

The recent advancement of large language models (LLMs), such as ChatGPT and GPT-4 (OpenAI, 2023), have sparked significant interest in developing LLM-based autonomous agents (AutoGPT, 2022) for complex task execution (Qin et al., 2023; Schick et al., 2023). Recent studies have explored the construction of text-based web browsing environments and how to instruct large language model agents to perform web navigation (Nakano et al., 2021; Gur et al., 2023; Zhou et al., 2023; Lu et al.,

2023). The primary challenge in these works lies in managing complex and verbose HTML texts, and solutions include simplifying and structuring HTML (Nakano et al., 2021; Zhou et al., 2023; Gur et al., 2023; Deng et al., 2023).

However, existing approaches overlook a critical functionality of browsing: rendering HTML into visual webpages. Particularly, vision capability is crucial for utilizing tools such as web browsers, as rendered web pages are inherently designed with user experience (UX), emphasizing intuitive information and structured presentation. This design principle of rendering makes visual analysis more effective than mere HTML representation. At present, large multimodal models (LMMs), particularly GPT-4V(ision) (OpenAI, 2023) and Gemini (Team et al., 2023), demonstrate a remarkable ability to integrate intricate visual cues with textual information. Existing studies such as Pix2Struct (Lee et al., 2023) and WebArena (Zhou et al., 2023), have initiated explorations into using screenshots as inputs for decision-making in web navigation, yet these are preliminary and do not represent a deep exploration. Therefore, building multimodal web agents to leverage the environment rendered by browsers through screenshots, thus mimicking human web browsing behavior, is now a viable approach to enhance web navigation abilities.

We introduce WebVoyager (Figure 1), a multimodal web agent designed to autonomously accomplish web tasks online from start to finish, managing the entire process end-to-end without any intermediate human intervention. WebVoyager processes the user query by making observations from screenshots and textual content in interactive web elements, formulates a thought on what action to take (such as clicking, typing, or scrolling, etc.), and then executes that action on the websites. Inspired by Set-of-Mark Prompting (Yang et al., 2023a), we mark interactive web elements on screenshots (see Figure 2) to facilitate decision-

\* Work done during the internship at Tencent AI Lab.

<sup>1</sup>Our code and data will be released at <https://github.com/MinorJerry/WebVoyager>

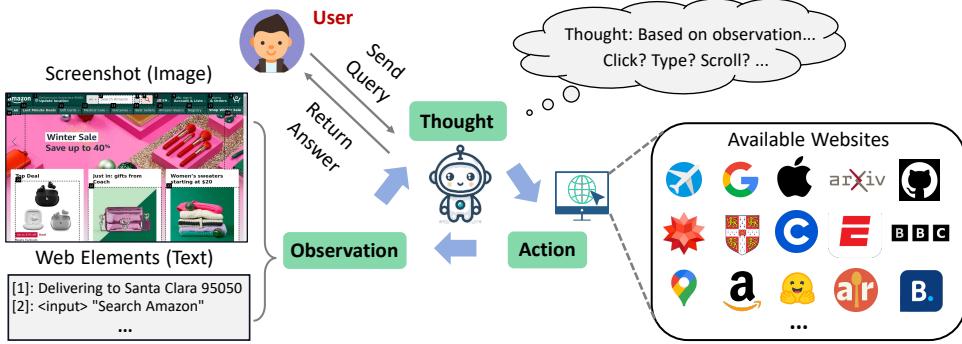


Figure 1: The overall workflow of WebVoyager. WebVoyager takes web tasks assigned by a human and automatically browses the web online. At each step, WebVoyager selects actions based on screenshots and text (the ‘type’ of the web element and its contents). Once the task is completed, the answers will be returned to the user. For example, for a user query: “Find the cost of a 2-year protection for PS4 on Amazon.”, the agent interacts with Amazon online, locates the PS4, identifies the 2-year protection price, and returns “\$30.99” to the user.

making for WebVoyager.

Another challenge is the evaluation of an end-to-end web agent. Existing benchmarks, such as Mind2Web (Deng et al., 2023), primarily focus on stepwise and offline evaluation, where agents follow a predefined “golden” trajectory for action selection. This approach, however, may not fully account for the variety of viable strategies to accomplish a task, as it only reflects one possible plan. This limitation could lead to a biased evaluation and difficulties in fairly comparing different methods. To accurately evaluate the capabilities of web agents in end-to-end task completion, we propose an automated evaluation protocol using GPT-4V. Specifically, we save screenshots throughout the online navigation process and then use GPT-4V to evaluate these trajectories together with the final results automatically. Human evaluations are also conducted to verify the results and the analysis shows that our evaluation protocol achieves 85.3% agreement with human judges, indicating GPT-4V can serve as a reliable evaluator for online agents.

We conduct evaluations on a newly collected dataset, which is semi-automatically generated using a self-instruct (Wang et al., 2022) method, comprising 643 web tasks from 15 commonly accessed websites. We also evaluate WebVoyager on 90 web-related tasks of level 1 and level 2 from the GAIA (Mialon et al., 2023), and 50 interactive open-web tasks from SeeAct (Zheng et al., 2024). We compare our WebVoyager with 1) GPT-4 (All Tools)<sup>2</sup>, and 2) WebVoyager in a text-only setting which employs the textual accessibility tree proposed in WebArena (Zhou et al., 2023) to describe web pages.

<sup>2</sup>GPT-4 (All Tools) is an integrated tool-based agent released by OpenAI in Oct. 2023. See <https://chat.openai.com/>

The results show that WebVoyager achieves a Task Success Rate of 59.1% on our new benchmark, significantly outperforming GPT-4 (All Tools) with a rate of 30.8% and the text-only setting with a rate of 40.1%, demonstrating the effectiveness of our method. Our research demonstrates the effectiveness of the WebVoyager method for web tasks, offering insights into the development of more intelligent and efficient web automation solutions.

## 2 Related Work

Autonomous web navigation (Shi et al., 2017; Yang et al., 2023b) requires an agent to follow instructions, construct plans, comprehend complex web structures, and decompose tasks into step-by-step decisions (Weng, 2023). To study web agents in a controlled environment, previous works constructed web simulators that contain simplified websites (Shi et al., 2017; Yao et al., 2022a). More recently, there has been a surge of interest in building more challenging and realistic benchmarks such as Mind2Web (Deng et al., 2023) and WebArena (Zhou et al., 2023).

Along with these new benchmarks, numerous efforts have been made to build autonomous web agents. WebGPT (Nakano et al., 2021) constructs a text-based web browsing environment and fine-tunes GPT-3 as a web agent. WebAgent (Gur et al., 2023) pretrains a T5 model to extract HTML snippets and leverages Flan-U-PaLM (Chowdhery et al., 2023) to generate Python code to interact with the environment. Besides fine-tuning, another line of work tries to build web agents by prompting LLMs (Yao et al., 2022b; Shinn et al., 2023; Ma et al., 2023). Multimodal web agents that integrate visual

signals have also been explored, WebGUM (Furuta et al., 2023) combines T5 (Raffel et al., 2020) with a Vision Transformer (ViT) to navigate using both screenshots and HTML text. PIX2ACT (Shaw et al., 2023) instead solely relies on web screenshots as inputs to predict agent actions. Unlike previous works that only consider a single modality or simplified web environments, we build a multi-modal agent that can complete tasks on real-world websites in this work. Concurrently with our work, SeeAct (Zheng et al., 2024) also leverages Large Multimodal Models (LMMs) for integrated visual understanding and actions on websites. However, the best SeeAct agent still relies on a finetuned cross-encoder model to select candidate elements for interaction. In contrast, WebVoyager does not require any additional modules.

### 3 WebVoyager

We aim to build an agent that can browse the open web autonomously without human intervention to complete user instructions. Given an instruction, our WebVoyager first instantiates a web browser and then performs actions with visual (i.e., screenshots) and textual (i.e., HTML elements) signals from the web. The agent produces an action based on the inputs at every step, which is then executed in the browser environment. The process continues until the agent decides to stop. The details of WebVoyager, including environment, interaction cycle, observation space, and action space, are as follows.

#### 3.1 Browsing Environment

We develop an automated web-browsing environment using Selenium<sup>3</sup>. Unlike WebArena (Zhou et al., 2023), we do not host any websites locally and allow the agent to explore the open web instead, which poses unique challenges such as floating ads, pop-up windows, constant updates, etc.<sup>4</sup> Still, we opt for online interaction with real websites as we believe that this setting truly reflects the real-world use cases (e.g., the agent needs access to real-time information from the web), and a successful web agent should be able to adapt to these challenges and consistently solve the problem robustly.

<sup>3</sup><https://www.selenium.dev/>

<sup>4</sup>Regarding CAPTCHAs (Completely Automated Public Turing test to tell Computers and Humans Apart) challenges, we believe it is important to respect the rules of these websites and prompt the agent to retrieve information from alternative sources.

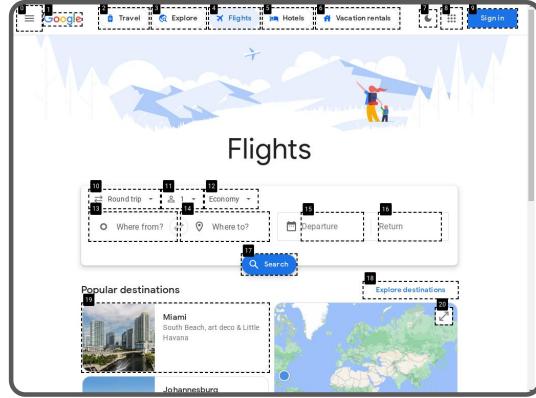


Figure 2: Examples of webpage screenshots provided to the agent. We add borders to most of the interactive elements on the web pages and label them with numerical tags in the top left corner.

#### 3.2 Interaction Formulation

Formally, we denote the Environment as  $\mathcal{E}$ , the large Multimodal Model as  $\mathcal{M}$ , the Observation Space as  $\mathcal{O}$ , and the Action Space as  $\mathcal{A}$ . At time step  $t$ , the model receives the context  $c_t$  as inputs, which consist of historical actions  $a_i$  and observations  $o_i$ , defined as:  $c_t = (o_1, a_1, \dots, o_{t-1}, a_{t-1}, o_t, I)$ . The model produces the action  $a_t$  at time  $t$ ,  $a_t = \mathcal{M}(c_t)$ , which is then executed in the environment. After execution, the environment sends back the observation at time  $t+1$ ,  $o_{t+1} = \mathcal{E}(o_t, a_t)$ . Then the context will be updated and this interaction process continues until the model generates a terminating action or the maximum step is reached.

Inspired by the paradigm of ReAct Prompting (Yao et al., 2022b), we also prompt our agent to generate a thought process first before generating the action code. Hence  $a_t$  can be further composed into  $(s_t, \hat{a}_t)$  where  $s_t$  and  $\hat{a}_t$  represent the natural language thought and action code respectively. Figure 7 in Appendix A presents the System Prompt we designed for the action prediction step. Also, it's worth noting that excessive observations of web pages from longer episodes may confuse the agent. Therefore, we perform context clipping to remove outdated web page information and only keep the three most recent observations in the inputs, and we keep the entire history of thoughts and actions to better guide the agent.

#### 3.3 Observation Space

Similar to how humans browse the web, our agent also takes the visual information from the web (screenshots) as the primary source of input. Using

screenshots allows us to avoid the burden of processing HTML DOM tree or accessibility tree to portray the overall structure of webpages, which can lead to overly verbose texts and impact the decision-making process of the agent. Inspired by Set-of-Mark Prompting (Yang et al., 2023a), we overlay bounding boxes of the interactive elements on the websites to better guide the agent’s action prediction. Unlike Yang et al. (2023a), we do not need any object detection module (Zou et al., 2023). Instead, we utilize GPT-4V-ACT<sup>5</sup>, a Javascript tool to extracts the interactive elements based on web element types and then overlays bounding boxes with numerical labels on the respective regions of the elements. GPT-4V-Act is efficient since it is rule-based without incorporating any object detection model.

As illustrated in Figure 2, the nature of webpages allows us to locate and outline each interactive element using this tool precisely. The numerical labels assigned to each element are also essential for the model to identify the elements requiring interaction, thereby facilitating accurate action determination. We empirically choose black color for the borders and the background of the labels to enhance clarity. We observe that using a single black color yields higher success rates than using multiple colors. We also provide the agent with auxiliary text as inputs, including the textual content embedded within the interactive element, the type of the element, and possibly some comment text in the aria-label attribute. To simplify the observation, we have disabled multiple tabs, i.e., all interactions occur within the current tab instead of opening new ones.

At every step, the agent receives the current screenshot, auxiliary text, and history as inputs, as discussed in (§3.2). In case the agent’s action raised an exception during execution, we additionally incorporated the error messages in the prompt and asked the model to regenerate the response. Note that each error correction attempt also consumes one step from the total exploration budget.

### 3.4 Action Space

We define the action space of our agent similar to how humans browse the web. To this end, we implement the most commonly used mouse and keyboard actions, sufficient for the agent to browse various web pages and locate the content required

for the task. With the help of numerical labels in screenshots, we enable the agent to respond with a concise Action Format. This approach precisely locates the elements requiring interaction and executes the corresponding actions. The usage of actions is as follows (more details in Appendix C): 1) **Click**. This action involves clicking on an element within a webpage, typically a link or a button. 2) **Input**. This composite action involves selecting a text box, deleting any existing content within it, and then inputting new content. 3) **Scroll**. Scrolling is a common operation for browsing webpages, usually involving the vertical movement of the entire page. 4) **Wait**. Action execution requires time, and this action is often used to wait for web pages to load. 5) **Back**. This action is used to return to the previous page. 6) **Jump to Search Engine**. There are often situations where agents get stuck at a certain website without finding an answer. This action enables the agent to jump to a search engine and start anew. 7) **Answer**. Once all questions in the task are resolved, this action concludes the iteration and provides an answer in line with the task requirements.

## 4 Benchmark for WebVoyager

### 4.1 Website Selection

We select 15 representative websites that cover different aspects of our daily life to ensure diversity in our evaluation, including Allrecipes, Amazon, Apple, ArXiv, BBC News, Booking, Cambridge Dictionary, Coursera, ESPN, GitHub, Google Flights, Google Map, Google Search, Huggingface, and Wolfram Alpha. Due to technical limitations, we regrettfully omit websites requiring login or CAPTCHA to access their content. Additionally, Google Search is a universal website that can serve as a starting point for any website, making our framework applicable to various scenarios.

### 4.2 Data Construction

We employ a combination of self-instruct (Wang et al., 2022) and human verification to construct our evaluation set. Figure 3 illustrates our data creation process. Initially, we manually sample and rewrite some tasks from Mind2Web (Yin et al., 2023; Deng et al., 2023) for websites including Google Flights, Google Map, Google Search, Booking, and Wolfram Alpha. This process yields initial seed tasks in the Task Pool for subsequent generation. In step two, we sample tasks from Task Pool as in-context

<sup>5</sup><https://github.com/ddupont808/GPT-4V-Act>

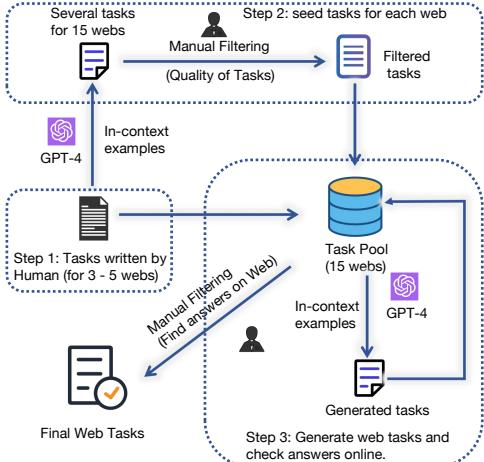


Figure 3: Data creation process using self-instruct.

examples (Dong et al., 2022) and prompt GPT-4 Turbo to generate approximately 100 new tasks (20 iterations). Then we manually verify each generated task and rewrite them if necessary to ensure its high quality and the answers can be found on the corresponding website, then we add them to the Task Pool as additional seed tasks. This step allows us to create human-validated seed tasks for each website. Finally, in step three, we sample more diverse in-context examples in the Task Pool and directly add the generated tasks to the Task Pool in each iteration. We manually verify that the generated tasks have low repetition, and the answers to the generated tasks can be found on the websites. In total, we collected 40+ tasks per website, resulting in a total of 643 tasks.

To further confirm that the generated tasks in the dataset have low repetition, We use the all-mpnet-base-v2<sup>6</sup> model to calculate pairwise similarity for 643 ques. Out of a total of 206,403 pairs, only 49 pairs have a similarity greater than 0.8, and 140 pairs have a similarity between 0.7 and 0.8. All of these have been manually checked and found to be acceptable. 99.68% of pairs have a similarity of less than 0.6. This demonstrates the diversity of our tasks and the robustness of our approach.

### 4.3 Annotation Process

After collecting the full task pool, we annotate answers for each task. Since some questions are open-ended and the web information may change, these questions may not have a fixed golden response. Thus, we label each data entry with an answer, categorized as “Possible” or “Golden.” For answers

<sup>6</sup><https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

labeled as “Golden,” we provide a comprehensive listing of possible responses and consider them stable in the short term. The “Possible” category covers the following scenarios: 1) Answers for open-ended tasks where it’s hard to find an exact match answer, such as summarization. 2) multiple answers satisfy the task, making it impractical to list all of them. Therefore, we provide a partial listing. 3) Tasks related to real-time information, where the answer might change, e.g., flight ticket prices. Hence, the “Possible” answers were also correct during our experiments. In total, 22.3% of questions are annotated with golden responses, and the rest only have possible answers.

## 5 Experiment

**Dataset and Metrics** Besides the benchmark introduced in Section 4, we also evaluated on 90 web browsing tasks (Level 1 and Level 2) from the GAIA dataset (Mialon et al., 2023), which also come with golden responses. Since GAIA does not provide specific websites, we instruct the agent to start with Google Search. We further evaluated our agent on the 50 tasks used in SeeAct agent’s online evaluation (Zheng et al., 2024) and compared with their results. Following WebArena (Zhou et al., 2023), the primary evaluation metric we adopt is the **Task Success Rate**, measuring the successful completion of tasks without considering whether the steps are optimal.

**Experimental Details** We employ GPT-4 Turbo with vision (gpt-4-vision-preview) as the backbone model of our agent, which showcases strong semantic and visual understanding capabilities equivalent to GPT-4V (OpenAI, 2023). Furthermore, we incorporate two additional backbone models, namely Claude 3 Opus (Anthropic, 2024) and GPT-4o (GPT-4 Omni, OpenAI, 2024), to augment the diversity of our research. We include the GPT-4 (All Tools) for baselines, which integrates vision, web browsing, code analysis, and various plugins in one model. Additionally, we consider a text-only baseline where the agent only receives the websites’ accessibility tree as input to prediction actions. We used a fixed browser window size of 1024 \* 768 pixels for our environment, ensuring a consistent size for the screenshots in our observations. We set the temperature to 1 during generation and allow the agent to explore at most 15 steps.

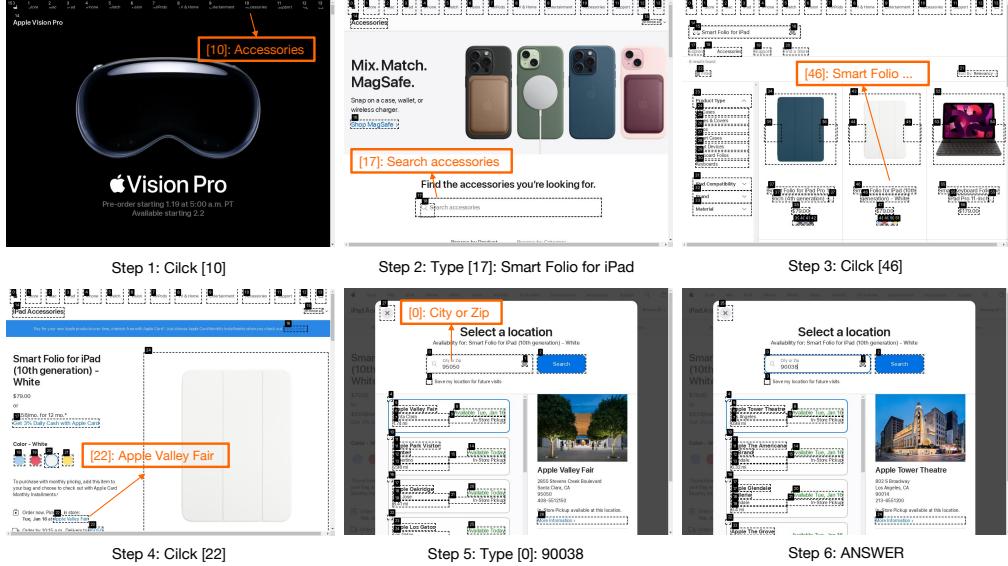


Figure 4: Screenshots of a complete trajectory of online web browsing. Given the task: ‘Search Apple for the accessory Smart Folio for iPad and check the closest pickup availability next to zip code 90038.’ The agent interacts with the Apple website and obtains the answer: ‘Apple Tower Theatre.’

## 5.1 Evaluation Methods

We adopt human evaluation as our main evaluation metric since most of the questions in our benchmark have open-ended answers. In particular, we provide the human evaluators with the complete trajectories of the agent’s interaction with the web (all screenshots and all actions), and ask them to provide a binary judgment of whether the agent successfully completed the task. For a subset of 300 tasks, we invite three annotators to judge each trajectory to understand the agreement among human annotators.

Even though human evaluations are accurate, they are often not scalable. Hence, we want to see if leveraging an LMM for automatic evaluation is feasible. To this end, we propose to use GPT-4V as an auto-evaluator that emulates the behavior of human evaluators to evaluate the navigation trajectories of WebVoyager. In particular, we provide the task, the responses from WebVoyager, and the last  $k$  screenshots to the evaluator and ask it to judge whether the agent has completed the task, where  $k$  is a hyper-parameter. The prompt of the GPT-4V evaluator is shown in Appendix B.

## 5.2 Result

Figure 4 presents an example that demonstrates how the agent interacts with the Apple website step by step in an online fashion to complete a task. In the final screenshot, the Agent acquires the desired information, then selects the “ANSWER” action to

respond and conclude the navigation. Additional examples are provided in the Appendix D.

We present the results for our dataset and the extracted GAIA web tasks in Table 1 and Figure 5. WebVoyager outperforms text-only and GPT-4 (All Tools) baselines by large margins in most website tasks, while it is slightly lower than Text-only on Allrecipes and similar to Text-only on Github, ESPN, Cambridge Dictionary and Wolfram Alpha. This is primarily because these websites are more text-heavy than others. Since WebVoyager mostly relies on web screenshots for decision-making, dense text might not be easily recognizable from the image. We think extracting such text from the HTML to augment the input could be a potential solution to this problem, suggesting a direction for future work. In Figure 5, WebVoyager also achieves much stronger performance than both baselines. Finally, WebVoyager has a success rate of 30% on the SeeAct online test set whereas the best SeeAct autonomous agent has 26%, showing the efficacy of our proposed agent.

We report the Agreement (the ratio of overlap) and Kappa ( $\kappa$ ; Cohen 1960) between consolidated human labels<sup>7</sup> and GPT-4V’s judgments on the subset of 300 tasks in Table 2. Here,  $k$  denotes the number of screenshots provided to GPT-4V, with “full” implying the full trajectory. GPT-4V’s agreement with human annotators gradually improves as

<sup>7</sup>the Fleiss’s Kappa (Fleiss, 1971) of human annotators before any discussion is 0.7, which is substantial agreement.

	Allrecipes	Amazon	Apple	ArXiv	GitHub	Booking	ESPN	Coursera
GPT-4 (All Tools)	11.1%	17.1%	44.2%	14.0%	48.8%	22.7%	31.8%	31.0%
WebVoyager <sub>Text-only</sub>	<b>55.6%</b>	31.7%	34.9%	32.6%	61.0%	2.3%	36.4%	23.8%
WebVoyager	53.3%	<b>58.5%</b>	<b>65.1%</b>	<b>51.2%</b>	<b>63.4%</b>	<b>43.2%</b>	<b>38.6%</b>	<b>73.8%</b>
<i>WebVoyager<sub>Text-only</sub>*</i>	57.8% $\pm$ 0.0%	43.1% $\pm$ 1.4%	36.4% $\pm$ 3.5%	50.4% $\pm$ 1.4%	63.4% $\pm$ 2.5%	2.3% $\pm$ 0.0%	38.6% $\pm$ 2.3%	24.6% $\pm$ 1.4%
<i>WebVoyager*</i>	51.1% $\pm$ 2.2%	52.9% $\pm$ 1.4%	62.8% $\pm$ 2.3%	52.0% $\pm$ 1.3%	59.3% $\pm$ 3.7%	32.6% $\pm$ 2.7%	47.0% $\pm$ 1.3%	57.9% $\pm$ 2.7%
<i>WebVoyager<sub>Claude</sub>*</i>	45.9% $\pm$ 3.4%	58.6% $\pm$ 4.2%	58.1% $\pm$ 4.0%	55.0% $\pm$ 7.0%	56.9% $\pm$ 1.4%	19.0% $\pm$ 1.3%	46.2% $\pm$ 1.3%	68.2% $\pm$ 1.3%
<i>WebVoyager<sub>GPT-4o</sub>*</i>	56.3% $\pm$ 1.3%	53.7% $\pm$ 2.5%	56.6% $\pm$ 1.3%	60.5% $\pm$ 0.0%	57.7% $\pm$ 3.7%	43.9% $\pm$ 3.5%	44.0% $\pm$ 2.7%	65.1% $\pm$ 2.8%
	Cambridge Dictionary	BBC News	Google Flights	Google Map	Google Search	Huggingface	Wolfram Alpha	Overall
GPT-4 (All Tools)	25.6%	9.5%	2.4%	53.7%	60.5%	37.2%	52.2%	30.8%
WebVoyager <sub>Text-only</sub>	62.8%	45.2%	7.1%	61.0%	67.4%	20.9%	58.7%	40.1%
WebVoyager	<b>65.1%</b>	<b>61.9%</b>	<b>59.5%</b>	<b>70.7%</b>	<b>76.7%</b>	<b>44.2%</b>	<b>63.0%</b>	<b>59.1%</b>
<i>WebVoyager<sub>Text-only</sub>*</i>	66.7% $\pm$ 3.6%	45.2% $\pm$ 2.4%	7.1% $\pm$ 0.0%	62.6% $\pm$ 2.8%	75.2% $\pm$ 1.3%	31.0% $\pm$ 1.4%	60.2% $\pm$ 1.3%	44.3% $\pm$ 0.6%
<i>WebVoyager*</i>	71.3% $\pm$ 1.3%	60.3% $\pm$ 2.8%	51.6% $\pm$ 1.4%	64.3% $\pm$ 2.8%	77.5% $\pm$ 2.7%	55.8% $\pm$ 2.3%	60.9% $\pm$ 2.2%	57.1% $\pm$ 0.2%
<i>WebVoyager<sub>Claude</sub>*</i>	71.3% $\pm$ 3.6%	66.7% $\pm$ 4.8%	15.1% $\pm$ 5.5%	55.3% $\pm$ 1.4%	72.9% $\pm$ 1.3%	53.5% $\pm$ 4.7%	51.5% $\pm$ 5.4%	52.8% $\pm$ 1.4%
<i>WebVoyager<sub>GPT-4o</sub>*</i>	82.2% $\pm$ 1.3%	54.8% $\pm$ 2.4%	28.6% $\pm$ 0.0%	56.9% $\pm$ 2.8%	63.6% $\pm$ 1.3%	42.6% $\pm$ 3.6%	65.2% $\pm$ 2.2%	55.5% $\pm$ 0.8%

Table 1: The main result for WebVoyager. Each website contains 40 to 45 tasks, and we report the Task Success Rate in the table. We show the results of GPT-4 (All Tools), WebVoyager<sub>Text-only</sub> using the accessibility tree, and WebVoyager by comparing with human expert labels. *WebVoyager\**, *WebVoyager<sub>Text-only</sub>\**, *WebVoyager<sub>Claude</sub>\** and *WebVoyager<sub>GPT-4o</sub>\** are results evaluated by GPT-4V (full trajectory, kappa = 0.70). For each automatic evaluation, we run GPT-4V evaluator three times to calculate the performance mean and standard deviation.

Success Rate	Consistency	
	Agreement	$\kappa$
k=1	47.7%	75.3%
k=2	55.3%	79.7%
k=3	54.3%	81.3%
Full	58.3%	85.3%
		0.51
		0.59
		0.62
		0.70

Table 2: Consistency between GPT-4V and Human. Success Rate is the overall success rate of all tasks given by GPT-4V. Based on the annotations given by GPT-4V and Human (after alignment), we compute Agreement, i.e., the label overlap, and the Kappa values.

it receives more information, and its final Kappa score also reaches 0.7, which is on par with the agreement among human annotators. The consistency between GPT-4V and humans suggests that GPT-4V is a promising automatic evaluator for multi-modal web agents. Accordingly, we report the automatic evaluation results based on GPT-4V in Table 1. The automatic evaluation results of three backbone models, GPT-4V, Claude 3 Opus, and GPT-4o, are relatively close, and their performance is significantly better than the Text-only setting (with GPT-4 as the backbone). However, there is a performance decline for both Claude and GPT-4o on Google Flights. Upon reviewing the trajectories, it is observed that GPT-4o falls into a cognitive bias, where it fails to correctly select the ‘one way’ option for one-way trip tasks, mistakenly assuming that only the departure date needs to be entered. On the other hand, Claude-3-Opus

WebVoyager Backbone	Evaluator		
	GPT-4V	Claude-3-Opus	GPT-4o
GPT-4V	57.1	55.1	63.0
Claude-3-Opus	52.8	61.6	55.4
GPT-4o	55.5	54.9	64.1

Table 3: Overall Task Success Rate of WebVoyager using automatic evaluation. We employ GPT-4V, Claude-3-Opus, and GPT-4o as backbones for WebVoyager and run all tasks, followed by automatic evaluation using these three models.

encounters difficulties in correctly interacting with web elements while inputting basic flight information. Modifying the system prompt for GPT-4o or Claude may potentially improve the performance.

Besides, we conduct the Claude-3-Opus based evaluation and the GPT-4o based evaluation. When provided with the full trajectory, the Claude-3-Opus achieves a kappa value of 0.6 with humans, indicating that it is less reliable than the GPT-4V. And the kappa value between GPT-4o and humans is 0.72, slightly higher than that of GPT-4V. Table 3 illustrates the Task Success Rate when using GPT-4V, Claude-3-Opus, and GPT-4o as backbones and auto evaluators. We observe that GPT-4o exhibits a more lenient attitude towards task performance results, while GPT-4V tends to be relatively strict. However, both models agree that Claude-3-Opus performs the worst. Claude-3-Opus, on the other hand, demonstrates a clear preference for its own results, believing that GPT-4V and GPT-4o are simi-

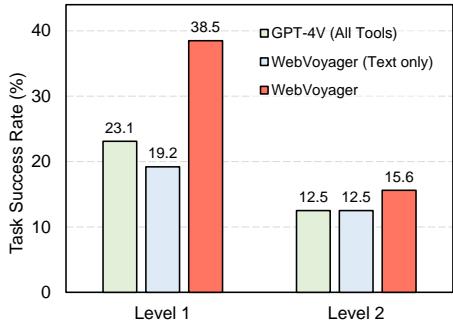


Figure 5: Success Rate for GAIA Level 1 and Level 2.

lar but considers itself to have the best performance. GPT-4o and GPT-4V also exhibit a certain bias towards their own results, with each considering itself to be superior to the other.

### 5.3 Discussions

#### Direct interaction with the websites is necessary

From our experience of using GPT-4 (All Tools), the primary limitation of GPT-4 (All Tools)’s performance is rooted in its reliance on Bing search for web browsing, predominantly depending on web pages fetched by Bing. It cannot directly access certain websites (such as Apple, Amazon, BBC News, etc.) for searching, clicking, or utilizing their sorting functions. This greatly limits the agent’s ability to complete certain types of tasks.

**Both text and vision are necessary for generalist web agents.** As discussed earlier, WebVoyager struggles with text-heavy websites. On the other hand, we observe that the text-only agent demonstrates significantly poorer performance on websites with complex visual elements, such as Booking and Flights, which require interactions with calendars and other intricate components. In these scenarios, the textual input such as the accessibility tree becomes highly complex and verbose, making it far less intuitive than using screenshots. Hence it’s necessary to incorporate both modalities of inputs when building the general purpose agents.

**Websites with more interactable elements are more challenging for agents.** We also calculate the average trajectory length of tasks completed within the maximum number of iterations, as well as the average number of interactive web elements present on the webpage screenshots. Figure 6 illustrates their relationship with the Task Success Rate. We posit that the average trajectory length serves as a measure of a task’s complexity to some extent, while the average number of numerical labels related to the decision-making process reflects

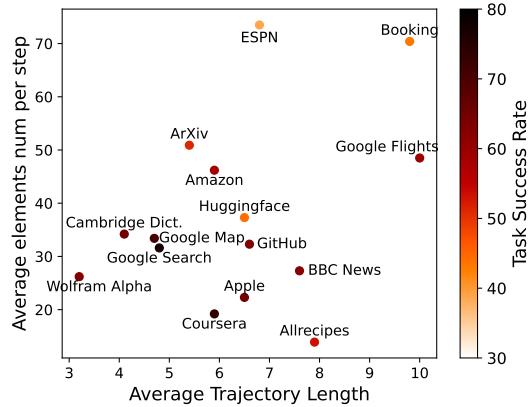


Figure 6: Factors related to task success rate. We show the average number of elements per page and the average trajectory length for each website type. The darker colors indicate a higher task success rate.

the complexity of a webpage. Intuitively, websites depicted in the lower-left corner of Figure 6, characterized by relatively simple webpages and shorter trajectory lengths, are expected to exhibit higher Task Success Rates. As observed in Figure 6, the results largely align with this intuition.

**Why not use Open Source models.** There are a few critical limitations of the existing open-sourced LMMs that prevent us from exploring other alternatives. Specifically, the web navigation task requires the agent to process fine-grained details from web page screenshots to make decisions, hence high-resolution is required to preserve the information from the web. However, most open-sourced LMMs such as LLaVA (Liu et al., 2024) reduce the image resolution to 224x224 or 336x336, which makes text with smaller fonts unrecognizable, hence they are unsuitable for web navigation tasks. Moreover, models like LLaVA have a max context length of 4096. In our case, the agent needs to handle trajectories as long as 15 steps, and it requires approximately 7000+ tokens, which does not fit in those models’ context size.

### 5.4 Error Analysis

In this section, we discuss and summarize the primary issues encountered by WebVoyager in the task completion process. These challenges will serve as critical entry points for future enhancements of the Task Success Rate and for devising strategies to obtain an Optimal Trajectory. We sampled 300 tasks from our benchmark and manually labeled the error category for each failed case, we show the distribution of errors in table 4. In Appendix F, we also provide specific examples for each issue.

Main reasons for Failure	Ratio
Navigation Stuck	44.4%
Visual Grounding Issue	24.8%
Hallucination	21.8%
Prompt Misalignment	9.0%

Table 4: Distribution of main failure reasons.

**Navigation Stuck** The most common failure is running out of steps before completing the task. There are three failure scenarios: 1) When the agent’s search query is not precise and explicit enough, it will be overwhelmed by irrelevant search results. The agent may prefer to browse different results or wait for incorrect outcomes rather than correct its previous action; 2) When the scrollable area is very small, the agent might not be able to locate the correct scrolling area and repeatedly request the execution of useless scrolling actions; 3) Sometimes in the middle of the page, the agent has trouble deciding whether to scroll up or down. The agent also tends to repeat its previous mistakes due to the input clipping, as mentioned in section 3.2. These meaningless or repetitive actions may hinder the completion of the task.

**Visual Grounding Issue** The visual grounding ability of our agent still has a large room for improvement. We observe the following issues: 1) The agent cannot interpret less frequently observed patterns, such as misidentifying characters representing the pronunciations or math formulas; 2) The agent doesn’t recognize the subtle difference between two observations and thinks the execution has failed; 3) The agent selects the wrong element for action execution due to proximity. For example, the model sometimes confuses adjacent elements and misinterprets numbers on a calendar as numerical labels. Sometimes textual information plays a significant role, offering valuable cues and assisting in distinguishing between overly dense web elements. We find that incorporating the text content included in Web Elements can alleviate these problems to some extent. However, a stronger visual encoder or additional text inputs might be needed.

**Hallucination** Agents sometimes produce seemingly correct answers, which may require checking carefully to identify errors. We mainly see the following two scenarios: 1) The agent may overlook certain task requirements and settle for an answer that is only partially correct. For instance, when asked for the cheapest product, the agent might

respond with a cheap product visible in a screenshot, neglecting the need to sort the options first. 2) The agent might execute a seemingly correct action without raising any errors, which deviate it from the correct reasoning path. For example, inputting content to the wrong text box when there are many text boxes on the webpage is still valid, yet it would guide the agent to obtain a wrong answer.

**Prompt Misalignment** Understanding and following complex prompts, as illustrated in Figure 7, often poses significant challenges. Moreover, longer trajectories may result in excessively lengthy contexts, hindering effective instruction following. Although many of the errors in Navigation Stuck and Hallucination categories can also be attributed to prompt design, we use Prompt Misalignment to categorize the following situations: 1) the agent fails to generate outputs that can be parsed into executable actions, e.g. providing only the ‘Thought’ without the corresponding ‘Action’; 2) Prematurely terminating the process using the ANSWER action, though agent knows that the task is not yet complete (explicitly mentioned in its answer).

## 6 Conclusion

We introduce WebVoyager, an innovative web agent powered by large multimodal models (LMM) that can complete real-world web tasks end-to-end by interacting with websites. We have shown through evaluations that WebVoyager outperforms several baselines by leveraging both visual and textual signals. We also propose an automatic evaluation protocol by leveraging GPT-4V as the evaluator for online agents. Our work demonstrates the promise of using advanced LMM capabilities in building intelligent web agents. We hope WebVoyager provides a strong foundation for future research toward building more versatile and capable web assistants.

## Limitations

We recognize the following limitations of our work. First, we haven’t supported all possible actions in our environment compared to actions a human user might take when browsing the web. e.g. the Drag action on web pages. Supporting such an action is challenging since the degree of Drag is not a finite set. We may allow it to choose the pixel values to be dragged if the Visual Grounding capabilities of LMMs are further enhanced. Second, our agent

currently can only analyze basic file formats (such as text files and PDF files) and doesn't support all file formats, especially videos. Enhancing support for additional file formats is a crucial step in the development of web agents and we leave it for future work.

Regarding the potential risks of our work, we believe that it requires a substantial amount of safety checks before deploying web agents like WebVoyager into real-world applications, as the agent might unintentionally download malicious content from unauthorized websites, or input private/confidential information on public websites. Also the agent might send fake requests to website servers or generate fake user activities, which might be harmful to website owners. Therefore it's necessary to take extra caution when using and testing our agent.

## Ethics Statement

Our experiments have been designed to operate within strict ethical guidelines. Specifically, we restrict our web agent to perform only non-login tasks. This approach is in full compliance with the terms of service and user agreements of the websites our agent interacts with. Furthermore, we closely monitor the agent's activities during its online evaluations. This monitoring is designed to identify and prevent any actions that could lead to potentially harmful consequences. By taking these precautions, we ensure that our research does not cross ethical boundaries or cause unintended harm.

Additionally, all task queries for evaluation undergo thorough manual inspection to ensure they are harmless and ethically sound. This manual inspection process is aimed at ensuring that every query is harmless and does not promote or propagate harmful content or actions. Our work aims to enhance user experience and accessibility while mitigating potential negative societal impacts. By proactively addressing ethics concerns, we dedicate to conducting research that benefits society while upholding high ethical standards.

## References

Armen Aghajanyan, Bernie Huang, Candace Ross, Vladimir Karpukhin, Hu Xu, Naman Goyal, Dmytro Okhonko, Mandar Joshi, Gargi Ghosh, Mike Lewis, et al. 2022. Cm3: A causal masked multi-modal model of the internet. *arXiv preprint arXiv:2201.07520*.

AI Anthropic. 2024. Introducing the next generation of claudie.

AutoGPT. 2022. [AutoGPT](#).

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.

Qi Chen, Dileepa Pitawela, Chongyang Zhao, Gengze Zhou, Hsiang-Ting Chen, and Qi Wu. 2023. [Webvln: Vision-and-language navigation on websites](#).

Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Yantao Li, Jianbing Zhang, and Zhiyong Wu. 2024. [Seeclick: Harnessing gui grounding for advanced visual gui agents](#).

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.

Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46.

Yong Dai, Duyu Tang, Liangxin Liu, Minghuan Tan, Cong Zhou, Jingquan Wang, Zhangyin Feng, Fan Zhang, Xueyu Hu, and Shuming Shi. 2022. One model, multiple modalities: A sparsely activated approach for text, sound, image, video and code. *arXiv preprint arXiv:2205.06126*.

Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Samuel Stevens, Boshi Wang, Huan Sun, and Yu Su. 2023. Mind2web: Towards a generalist agent for the web. *arXiv preprint arXiv:2306.06070*.

Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. 2022. A survey for in-context learning. *arXiv preprint arXiv:2301.00234*.

Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.

Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological bulletin*, 76(5):378.

- Hiroki Furuta, Ofir Nachum, Kuang-Huei Lee, Yutaka Matsuo, Shixiang Shane Gu, and Izzeddin Gur. 2023. Multimodal web navigation with instruction-finetuned foundation models. *arXiv preprint arXiv:2305.11854*.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913.
- Izzeddin Gur, Hiroki Furuta, Austin Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck, and Aleksandra Faust. 2023. A real-world webagent with planning, long context understanding, and program synthesis. *arXiv preprint arXiv:2307.12856*.
- Jack Hessel, Jena D Hwang, Jae Sung Park, Rowan Zellers, Chandra Bhagavatula, Anna Rohrbach, Kate Saenko, and Yejin Choi. 2022. The abduction of sherlock holmes: A dataset for visual abductive reasoning. In *European Conference on Computer Vision*, pages 558–575. Springer.
- Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Chong Lim, Po-Yu Huang, Graham Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and Daniel Fried. 2024. *Visualwebarena: Evaluating multimodal agents on realistic visual web tasks*.
- Kenton Lee, Mandar Joshi, Iulia Raluca Turc, Hexiang Hu, Fangyu Liu, Julian Martin Eisenschlos, Urvashi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova. 2023. Pix2struct: Screenshot parsing as pretraining for visual language understanding. In *International Conference on Machine Learning*, pages 18893–18912. PMLR.
- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. 2019. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024. Visual instruction tuning. *Advances in neural information processing systems*, 36.
- Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. *Advances in Neural Information Processing Systems*, 35:2507–2521.
- Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. 2023. Chameleon: Plug-and-play compositional reasoning with large language models. *arXiv preprint arXiv:2304.09842*.
- Kaixin Ma, Hongming Zhang, Hongwei Wang, Xiaoman Pan, and Dong Yu. 2023. Laser: Llm agent with state-space exploration for web navigation. *arXiv preprint arXiv:2309.08172*.
- Grégoire Mialon, Clémentine Fourrier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom. 2023. Gaia: a benchmark for general ai assistants. *arXiv preprint arXiv:2311.12983*.
- Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. 2021. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*.
- OpenAI. 2023. *Gpt-4 technical report*.
- OpenAI. 2024. Hello gpt-4o.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. 2023. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*.
- Peter Shaw, Mandar Joshi, James Cohan, Jonathan Berant, Panupong Pasupat, Hexiang Hu, Urvashi Khandelwal, Kenton Lee, and Kristina Toutanova. 2023. From pixels to ui actions: Learning to follow instructions via graphical user interfaces. *arXiv preprint arXiv:2306.00245*.
- Tianlin Shi, Andrej Karpathy, Linxi Fan, Jonathan Hernandez, and Percy Liang. 2017. World of bits: An open-domain platform for web-based agents. In *International Conference on Machine Learning*, pages 3135–3144. PMLR.
- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik R Narasimhan, and Shunyu Yao. 2023. Reflexion: Language agents with verbal reinforcement learning. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.

- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hanneh Hajishirzi. 2022. Self-instruct: Aligning language model with self generated instructions. *arXiv preprint arXiv:2212.10560*.
- Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. 2021. Simvlm: Simple visual language model pretraining with weak supervision. *arXiv preprint arXiv:2108.10904*.
- Lilian Weng. 2023. [Llm-powered autonomous agents.](https://lilianweng.github.io) [lilianweng.github.io](https://lilianweng.github.io).
- An Yan, Zhengyuan Yang, Wanrong Zhu, Kevin Lin, Linjie Li, Jianfeng Wang, Jianwei Yang, Yiwu Zhong, Julian McAuley, Jianfeng Gao, Zicheng Liu, and Lijuan Wang. 2023. Gpt-4v in wonderland: Large multimodal models for zero-shot smartphone gui navigation.
- Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chunyu Li, and Jianfeng Gao. 2023a. Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v. *arXiv preprint arXiv:2310.11441*.
- Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. 2023b. The dawn of Imms: Preliminary explorations with gpt-4v (ision). *arXiv preprint arXiv:2309.17421*, 9(1).
- Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022a. Webshop: Towards scalable real-world web interaction with grounded language agents.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022b. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*.
- Da Yin, Faeze Brahman, Abhilasha Ravichander, Khyathi Chandu, Kai-Wei Chang, Yejin Choi, and Bill Yuchen Lin. 2023. Lumos: Learning agents with unified data, modular design, and open-source llms. *arXiv preprint arXiv:2311.05657*.
- Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. From recognition to cognition: Visual commonsense reasoning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6720–6731.
- Chi Zhang, Zhao Yang, Jiaxuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. 2023. Appagent: Multimodal agents as smartphone users.
- Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. 2024. Gpt-4v (ision) is a generalist web agent, if grounded. *arXiv preprint arXiv:2401.01614*.
- Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Yonatan Bisk, Daniel Fried, Uri Alon, et al. 2023. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854*.
- Zhengxia Zou, Keyan Chen, Zhenwei Shi, Yuhong Guo, and Jieping Ye. 2023. Object detection in 20 years: A survey. *Proceedings of the IEEE*.

## A Prompt for WebVoyager

The System Prompt for WebVoyager is shown in Figure 7. The Prompt’s guidelines hold potential for optimization and should be generalist rather than website-specific in design. Incorporating specific issues from websites into the system prompt may compromise the agent’s universality.

## B Prompt for Auto Evaluation

Figure 8 demonstrates using GPT-4V as an evaluator for web tasks, involving web task instruction, screenshots in the trajectory, and WebVoyager responses. We require GPT-4V to mark success or not success. The temperature is set to 0 to reduce randomness during evaluation.

## C Action Space

We detail the interaction actions that WebVoyager employs to navigate and operate within web environments. These actions are fundamental to how the agent interacts with web pages, retrieves information, and performs specific tasks as part of its operational protocol. The actions range from basic web navigation to more complex operations, tailored to efficiently gather data and respond to queries. Each action is designed with a specific format for easy identification and execution.

- Click. This action involves clicking on an element within a webpage, typically a link or a button. If clicking a link results in the download of a PDF file, we automatically parse its content using the OpenAI Assistant API<sup>8</sup> and incorporate it into the Observation. Action Format: Click [Numerical\_Label].
- Input. This is a composite action that involves selecting a text box, deleting any existing content within it, and then inputting new content. To minimize interaction frequency, an automatic ENTER key press follows the input completion. Action Format: Type [Numerical\_Label]; [Content].
- Scroll. Scrolling is a common operation for browsing webpages, usually involving the vertical movement of the entire page. However, there are instances where only a specific section within the webpage is scrollable. In such

cases, we expect the Agent to select an element within the scrollable area and navigate to that particular region for scrolling. Action Format: Scroll [Numerical\_Label or WINDOW]; [up or down].

- Wait. Action execution requires time, and this action is often used to wait for web pages to load. Action Format: Wait.
- Back. This action is used to return to the previous page. We consider the forward action unnecessary because it can be achieved by repeating previous actions. Action Format: GoBack.
- Jump to Search Engine. There are often situations where agents get stuck at a certain website, without finding an answer. This action enables the agent to jump to a search engine and start anew. In this work, we just adopt Google Search. Action Format: Google.
- Answer. Once all questions in the task are resolved, this action concludes the iteration and provides an answer in line with the task requirements. Action Format: ANSWER; [Content].

## D Additional Trajectories

In Figure 4, we demonstrate how WebVoyager automatically navigates the Apple website and accomplishes the task. In this section, we exhibit the trajectories for the remaining websites that complete web tasks successfully. We provide a screenshot for each step, accompanied by the action generated by WebVoyager. The specific navigation trajectories for each website are illustrated in Figures 9 to 22. In addition, we also explore the performance of WebVoyager on websites in other languages, and we provide two example trajectories in Chinese and Spanish in Figure 23 and Figure 24.

## E Additional Related Work

**Vision-based Agents** Concurrent to our work, a few related works also studied vision-based autonomous agents. VisualWebArena (Koh et al., 2024) extends WebArena with additional websites and tasks that focus on visual reasoning to facilitate research on vision-based web agents. SeeClick (Cheng et al., 2024) focused on finetuning an LMM to solely leverage screenshots as inputs to interact

<sup>8</sup><https://platform.openai.com/docs/assistants/overview>

Imagine you are a robot browsing the web, just like humans. Now you need to complete a task. In each iteration, you will receive an Observation that includes a screenshot of a webpage and some texts. This screenshot will feature Numerical Labels placed in the TOP LEFT corner of each Web Element. Carefully analyze the visual information to identify the Numerical Label corresponding to the Web Element that requires interaction, then follow the guidelines and choose one of the following actions:

1. Click a Web Element.
2. Delete existing content in a textbox and then type content.
3. Scroll up or down.

...

Correspondingly, Action should STRICTLY follow the format:

- Click [Numerical\_Label]
- Type [Numerical\_Label]; [Content]
- Scroll [Numerical\_Label or WINDOW]; [up or down]
- Wait
- GoBack
- Google
- ANSWER; [content]

Key Guidelines You MUST follow:

\* Action guidelines \*

- 1) Execute only one action per iteration.

...

\* Web Browsing Guidelines \*

- 1) Don't interact with useless web elements like Login, Sign-in, donation that appear in Webpages.

...

Your reply should strictly follow the format:

Thought: {Your brief thoughts (briefly summarize the info that will help ANSWER)}

Action: {One Action format you choose}

Then the User will provide:

Observation: {A labeled screenshot Given by User}

Figure 7: System Prompt for WebVoyager. We instruct agents to perform web navigation, along with specific browsing actions and action formats. To enhance efficiency and accuracy, we can incorporate additional general into the prompts. These guidelines should be generic and not about a specific website to ensure generalizability.

with websites. WebVNL (Chen et al., 2023) introduced a web simulator that provides both HTML text and screenshots to finetune supervised vision-language models. GPT-4V Wonderland (Yan et al., 2023) and AppAgent (Zhang et al., 2023) instead focus on building agents that can operate smartphone apps using the GPT-4V as the backbone. These works further underscore the promising prospects in this field.

**Large Multimodal Models.** In recent years, significant strides have been made in unifying image and text representations within a single multimodal model through joint training with image and text (Li et al., 2019; Dosovitskiy et al., 2020; Wang et al., 2021; Dai et al., 2022; Aghajanyan et al., 2022). Large Multimodal Models (LMMs), following in the footsteps of Large Language Models (Brown et al., 2020; Chen et al., 2021; Chowdhery et al., 2023), attain the capability of instruction following (Ouyang et al., 2022) and exhibit robust multimodal comprehension. Represented by GPT-4V (OpenAI, 2023) and Gemini (Team et al., 2023),

LMMs have demonstrated impressive performance on benchmarks (Goyal et al., 2017; Lu et al., 2022; Zellers et al., 2019; Hessel et al., 2022), establishing a foundation for the construction of multimodal agents in subsequent research.

## F Error Cases

In this section, we provide specific examples of the four types of errors mentioned in the Error Analysis section. Figure 25 illustrates a snippet of WebVoyager navigating on Google Flights and a Visual Grounding Issue appears. The task is to retrieve one-way flight information for January 22nd; however, it selects December 22nd on the Calendar and fails to make the necessary corrections. Although it attempts to modify the date in step 6, it ultimately fails to do so. Figure 26 illustrates a situation of WebVoyager navigating on Allrecipes, encountering the Navigation Stuck issue. The agent requires multiple downward scrolls to locate the correct ingredients. However, it experiences confu-

As an evaluator, you will be presented with three primary components to assist you in your role:

1. Web Task Instruction: This is a clear and specific directive provided in natural language, detailing the online activity to be carried out. These requirements may include conducting searches, verifying information, comparing prices, checking availability, or any other action relevant to the specified web service (such as Amazon, Apple, ArXiv, BBC News, Booking etc).
  2. Result Screenshots: This is a visual representation of the screen showing the result or intermediate state of performing a web task. It serves as visual proof of the actions taken in response to the instruction.
  3. Result Response: This is a textual response obtained after the execution of the web task. It serves as textual result in response to the instruction.
- You DO NOT NEED to interact with web pages or perform actions such as booking flights or conducting searches on websites.  
-- You SHOULD NOT make assumptions based on information not presented in the screenshot when comparing it to the instructions.  
-- Your primary responsibility is to conduct a thorough assessment of the web task instruction against the outcome depicted in the screenshot and in the response, evaluating whether the actions taken align with the given instructions.  
-- NOTE that the instruction may involve more than one task, for example, locating the garage and summarizing the review. Failing to complete either task, such as not providing a summary, should be considered unsuccessful.  
-- NOTE that the screenshot is authentic, but the response provided by LLM is generated at the end of web browsing, and there may be discrepancies between the text and the screenshots.  
-- Note the difference: 1) Result response may contradict the screenshot, then the content of the screenshot prevails, 2) The content in the Result response is not mentioned on the screenshot, choose to believe the content.

You should elaborate on how you arrived at your final evaluation and then provide a definitive verdict on whether the task has been successfully accomplished, either as 'SUCCESS' or 'NOT SUCCESS'.

Figure 8: System Prompt for Auto Evaluation by GPT-4V.

sion during the process, and it is uncertain whether to scroll up or down. Figure 27 depicts the Hallucination issue encountered by WebVoyager on the Coursera website. In the task, we query the number of quizzes in the "Artificial Intelligence for Healthcare" course. However, the agent only identifies the quiz in module 1 of the course, which is not the optimal answer and does not fulfill the task requirements. Figure 28 illustrates the issue of Prompt Misalignment encountered while browsing BBC News. WebVoyager fails to fulfill all the task requirements. Instead of completing the navigation, it provides partial answers and tells me how to find complete answers, which is not end-to-end.

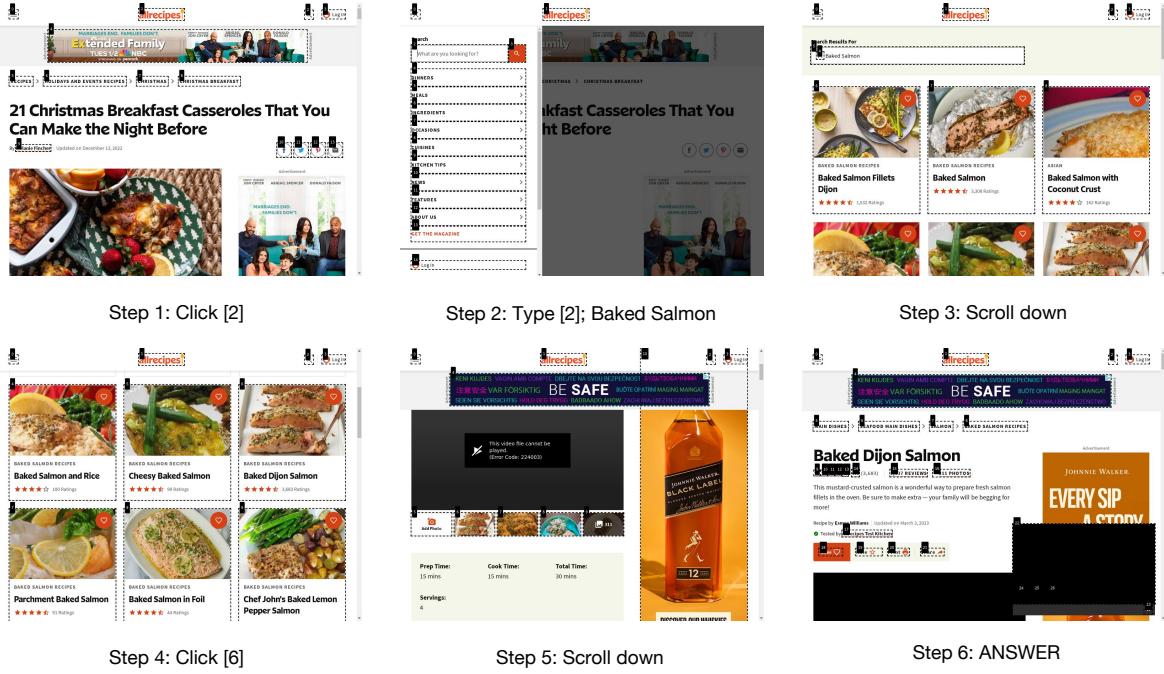


Figure 9: Screenshots of a complete trajectory of browsing Allrecipes. Given the task: “Find a recipe for Baked Salmon that takes less than 30 minutes to prepare and has at least a 4-star rating based on user reviews.” The agent interacts with the Allrecipes website and obtains the answer: “The “Baked Dijon Salmon” recipe meets the user’s criteria, with a 4.6-star rating and a preparation time of 15 minutes.”



Figure 10: Screenshots of a complete trajectory of browsing Amazon. Given the task: “Search for an Xbox Wireless controller with green color and rated above 4 stars.” The agent interacts with the Amazon website and obtains the answer: “The green Xbox Wireless controller (“Xbox Core Wireless Gaming Controller – Velocity Green”) rated above 4 stars has been found on Amazon with a rating of 4.7 out of 5 stars.”

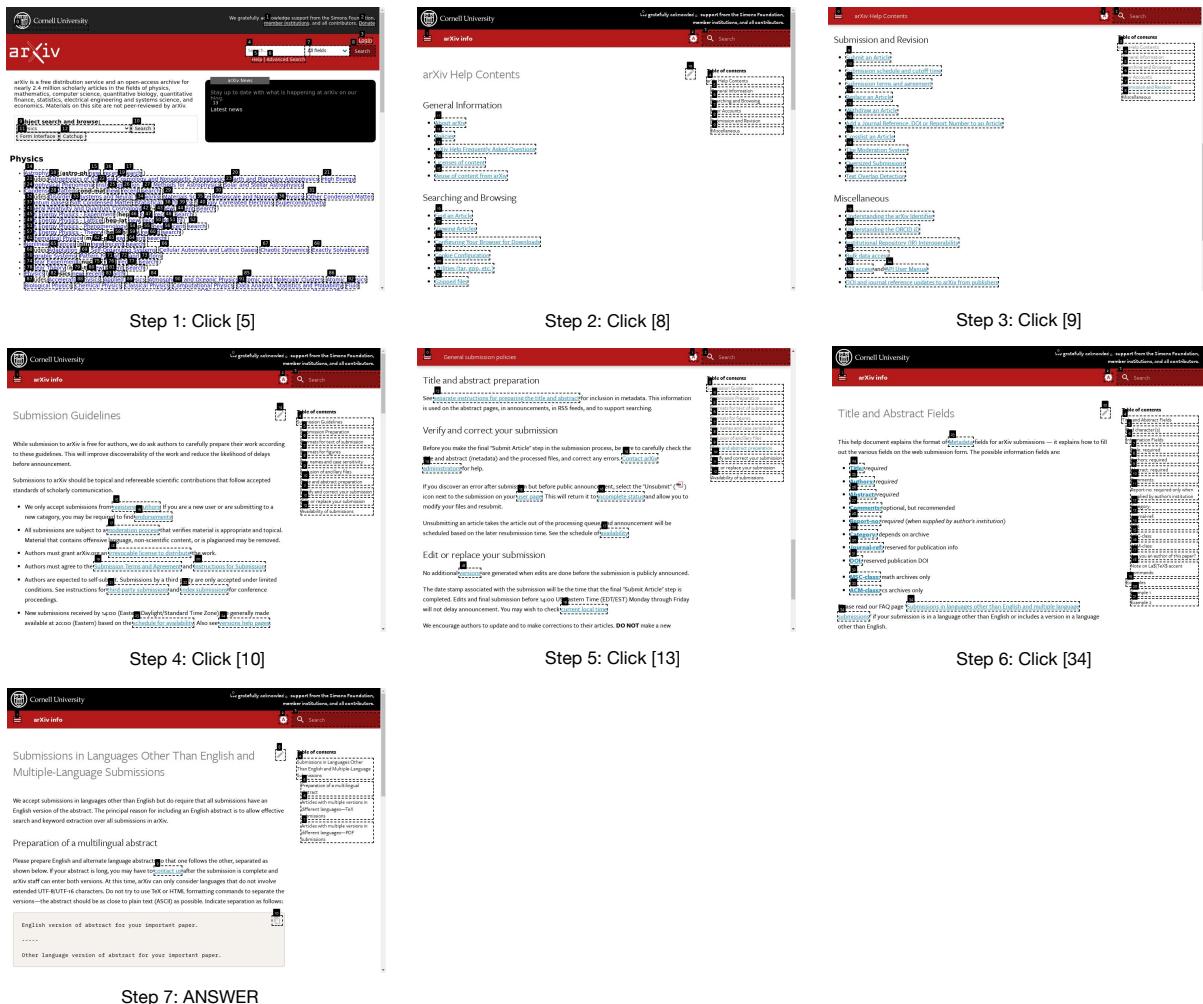


Figure 11: Screenshots of a complete trajectory of browsing ArXiv. Given the task: “For non-English submissions, do I need to provide a multi-language abstract, if needed, answer the separator between the multiple abstracts.” The agent interacts with the ArXiv website and obtains the answer: “Line (“-----”)”

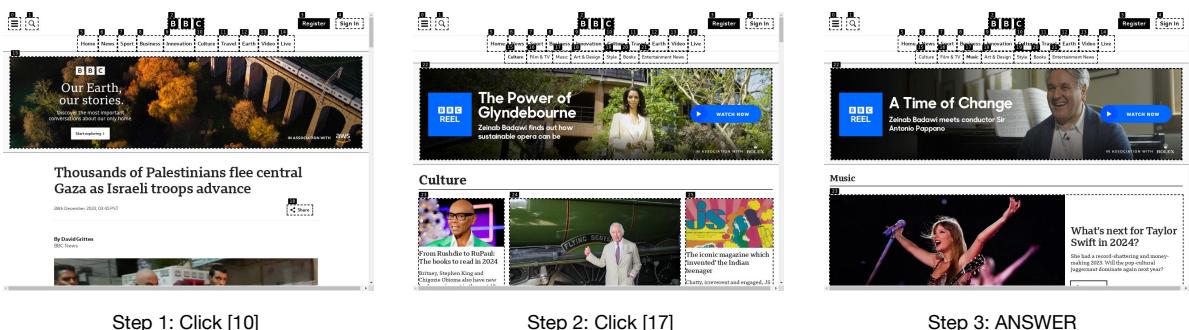


Figure 12: Screenshots of a complete trajectory of browsing BBC News. Given the task: “Find out which musician made the headlines in Music News.” The agent interacts with the BBC News website and obtains the answer: “The musician who made the headlines in Music News is Taylor Swift.”



Figure 13: Screenshots of a complete trajectory of browsing Booking. Given the task: “Find the cheapest available hotel room for a three-night stay from 1st Jan in Jakarta. The room is for 2 adults, just answer the cheapest hotel room and the price.” The agent interacts with the Booking website and obtains the answer: “The cheapest hotel room is at OYO 3755 Sweet Home for US\$14 for a three-night stay.”

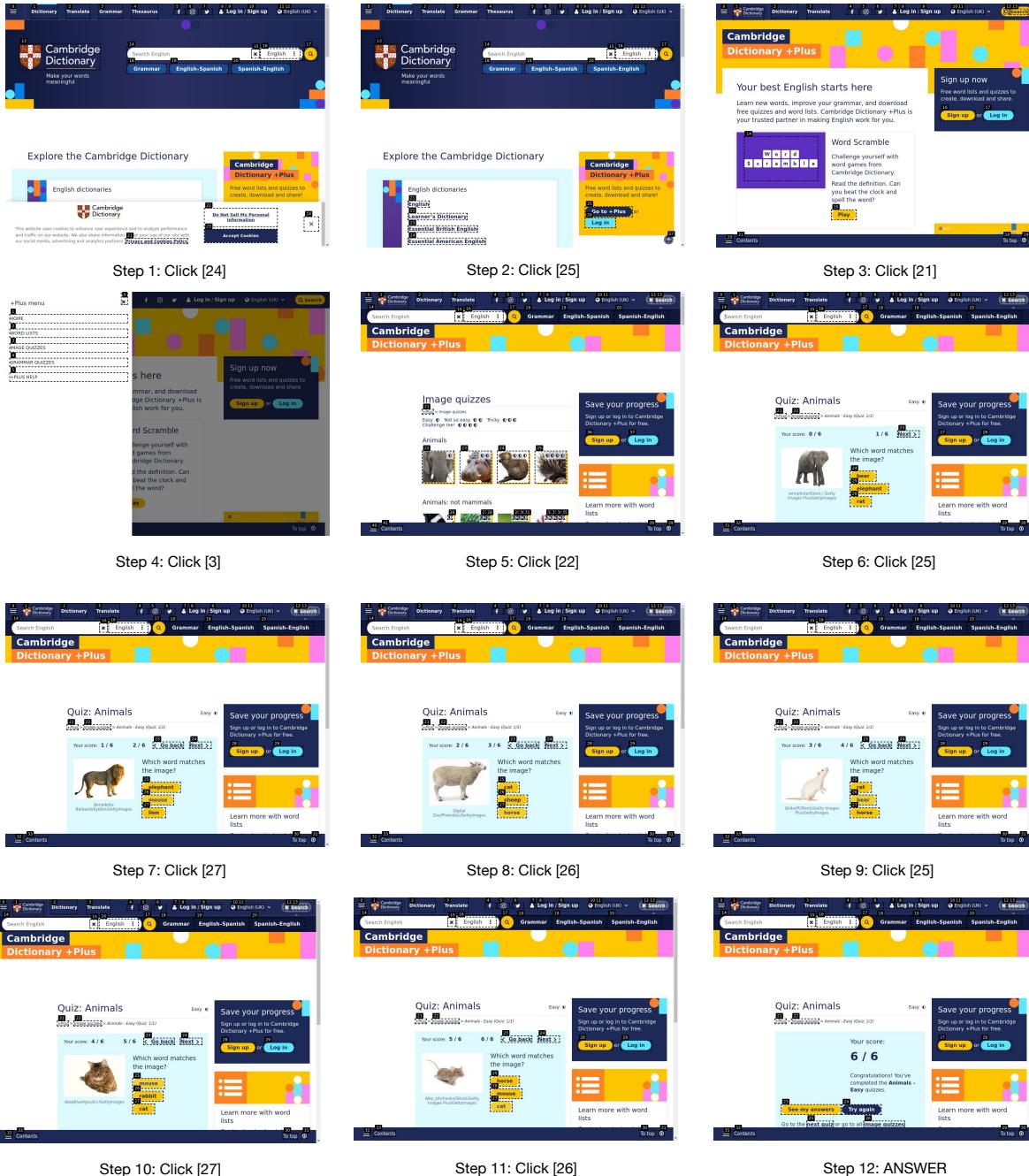


Figure 14: Screenshots of a complete trajectory of browsing Cambridge Dictionary. Given the task: “Go to the Plus section of Cambridge Dictionary, find Image quizzes and do an easy quiz about Animals and tell me your final score.” The agent interacts with the Cambridge Dictionary website and obtains the answer: “The final score on the easy Animals quiz is 6 out of 6.”

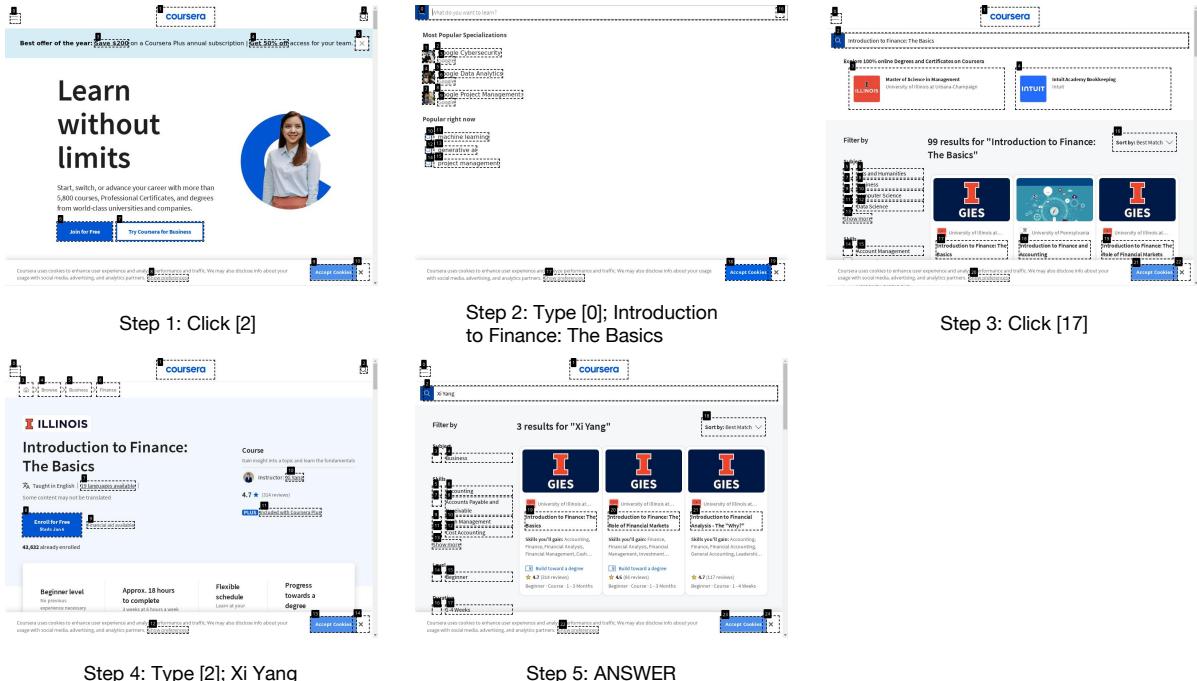


Figure 15: Screenshots of a complete trajectory of browsing Coursera. Given the task: “Identify a course on Coursera named ‘Introduction to Finance: The Basics’, who is the course instructor, and what other courses does he/she teach.” The agent interacts with the Coursera website and obtains the answer: The course instructor for ‘Introduction to Finance: The Basics’ is Xi Yang. Xi Yang also teaches ‘Introduction to Finance: The Role of Financial Markets’ and ‘Introduction to Financial Analysis - The “Why?”’”

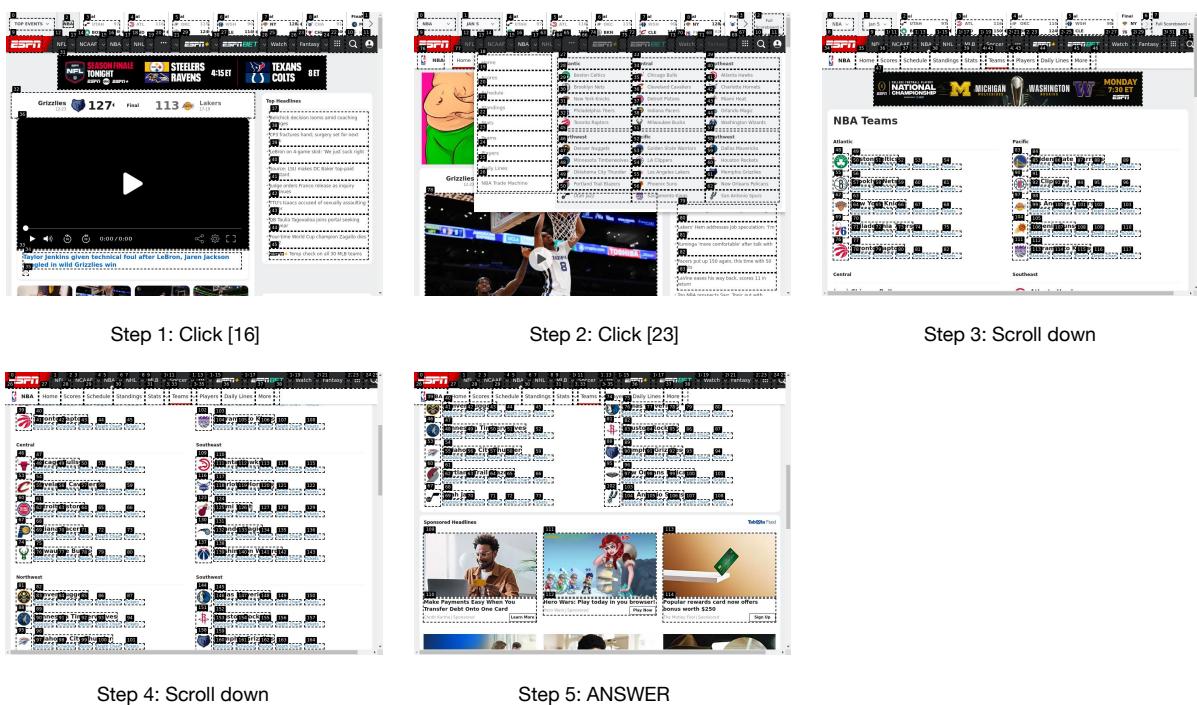


Figure 16: Screenshots of a complete trajectory of browsing ESPN. Given the task: “How many NBA teams are there and list all the teams with ‘New’ in their name.” The agent interacts with the ESPN website and obtains the answer: “There are 30 NBA teams. The teams with ‘New’ in their name are the New York Knicks and the New Orleans Pelicans.”

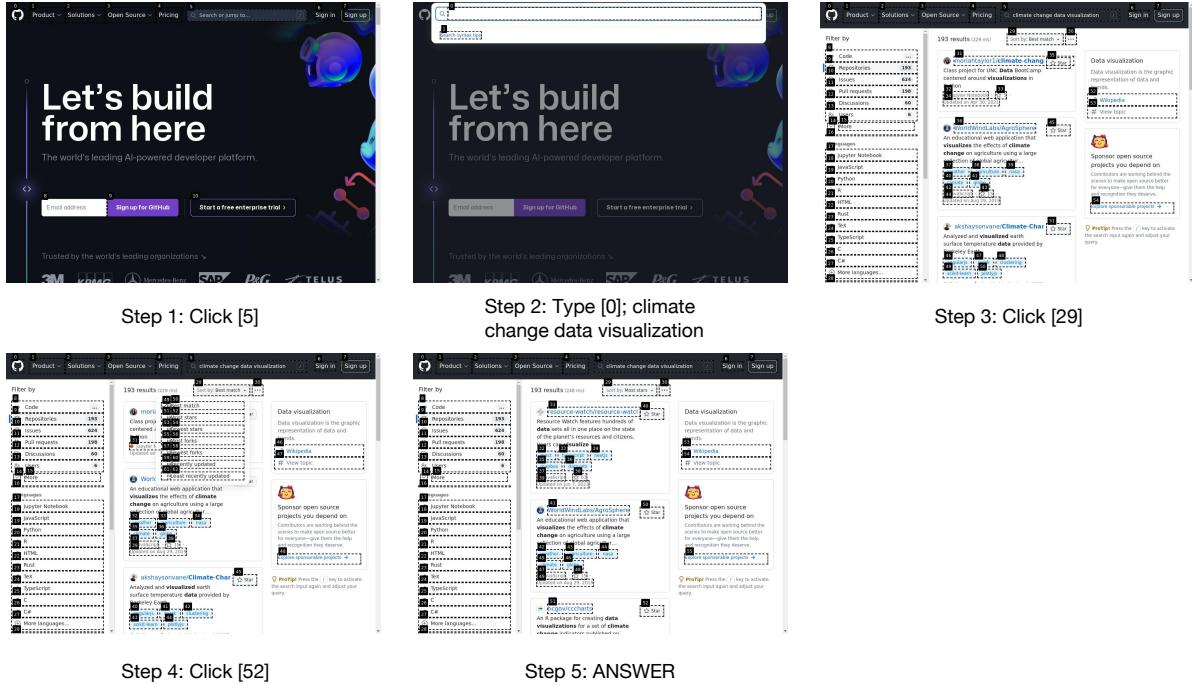


Figure 17: Screenshots of a complete trajectory of browsing GitHub. Given the task: “Search for an open-source project related to ‘climate change data visualization’ on GitHub and report the project with the most stars.” The agent interacts with the GitHub website and obtains the answer: “resource-watch/resource-watch with 63 stars”

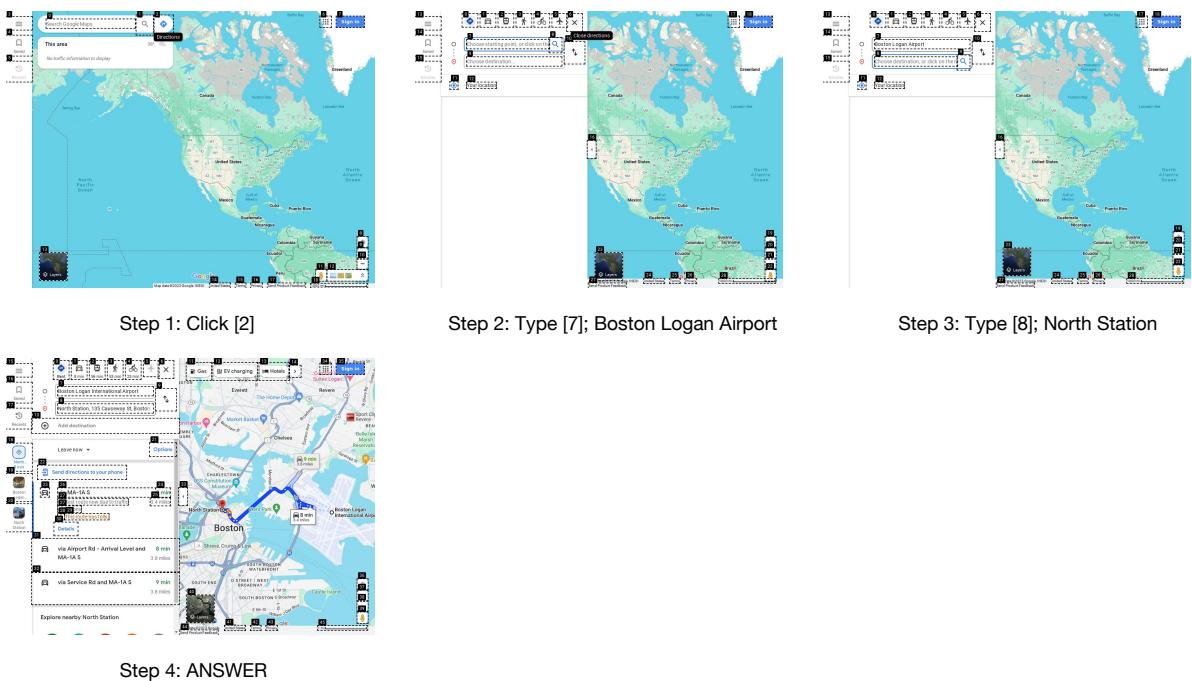


Figure 18: Screenshots of a complete trajectory of browsing Google Map. Given the task: “Plan a trip from Boston Logan Airport to North Station.” The agent interacts with the Google Map website and obtains the answer: “The fastest route from Boston Logan International Airport to North Station is via MA-1A S and takes approximately 8 minutes in current traffic conditions.”

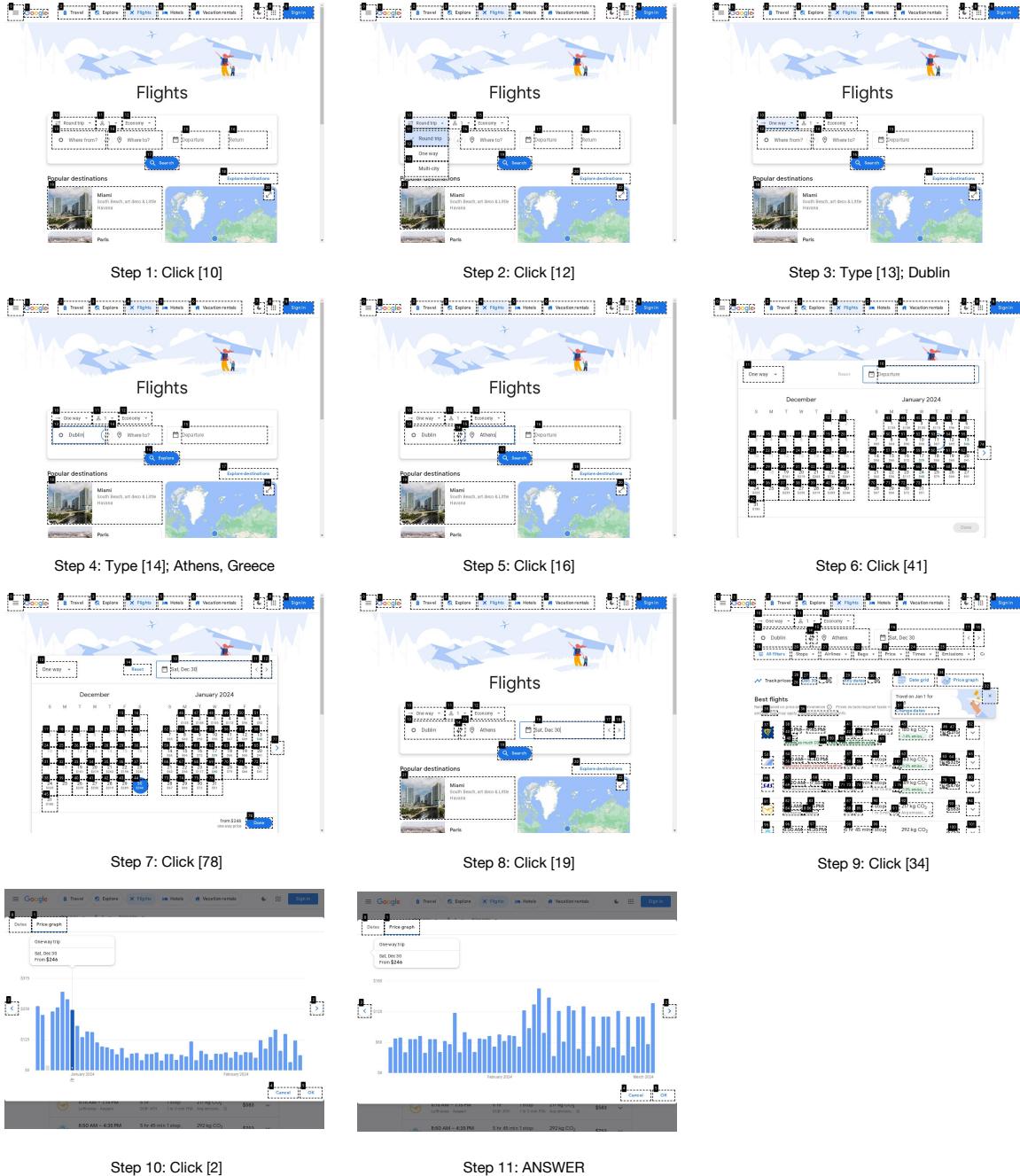


Figure 19: Screenshots of a complete trajectory of browsing Google Flights. Given the task: “Search a one-way flight from Dublin To Athens Greece for 1 Adult that leaves on December 30 and analyse the price graph for the next 2 months.” The agent interacts with the Google Flights website and obtains the answer: “The price graph has been analyzed, and it shows the trends for flight prices from Dublin to Athens over the next two months.”

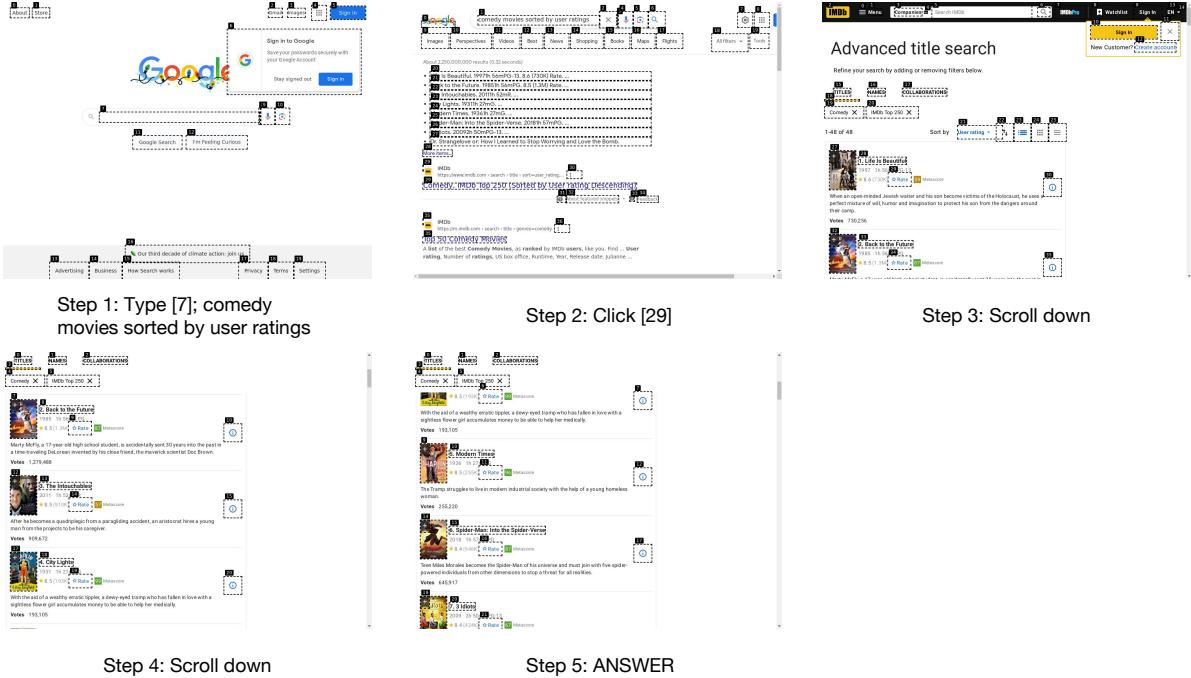


Figure 20: Screenshots of a complete trajectory of browsing Google Search. Given the task: “Show me a list of comedy movies, sorted by user ratings. Show me the Top 5 movies.” The agent interacts with the Google Search website and obtains the answer: “The Top 5 comedy movies sorted by user ratings are:Life Is Beautiful; Back to the Future; The Intouchables; City Lights; Modern Times”

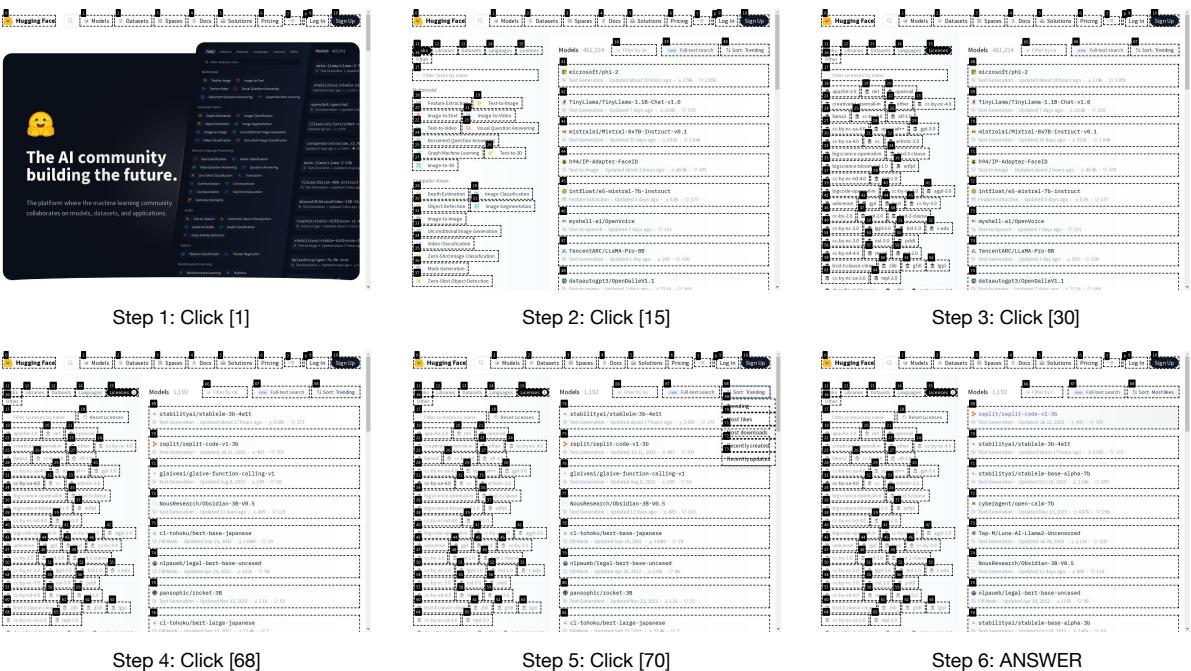


Figure 21: Screenshots of a complete trajectory of browsing Huggingface. Given the task: “Look up a model with a license of cc-by-sa-4.0 with the most likes on Huggingface.” The agent interacts with the Huggingface website and obtains the answer: “The model with a license of cc-by-sa-4.0 with the most likes on Hugging Face is “replit/replit-code-v1-3b” with 703 likes.”

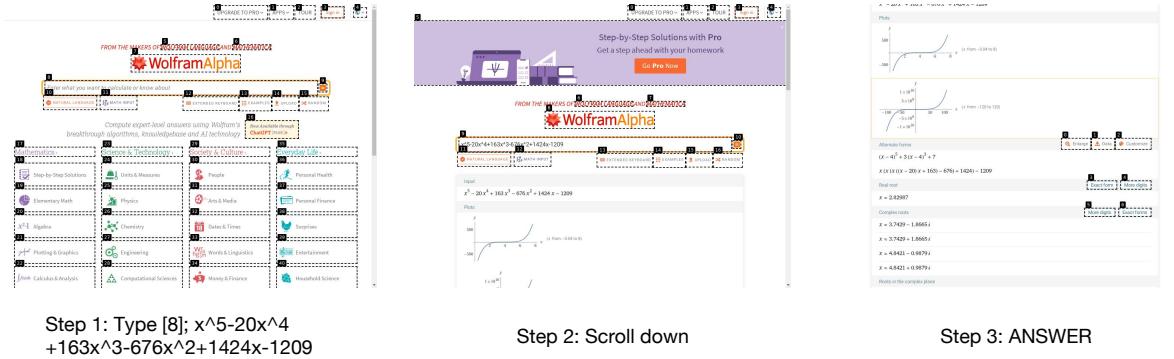


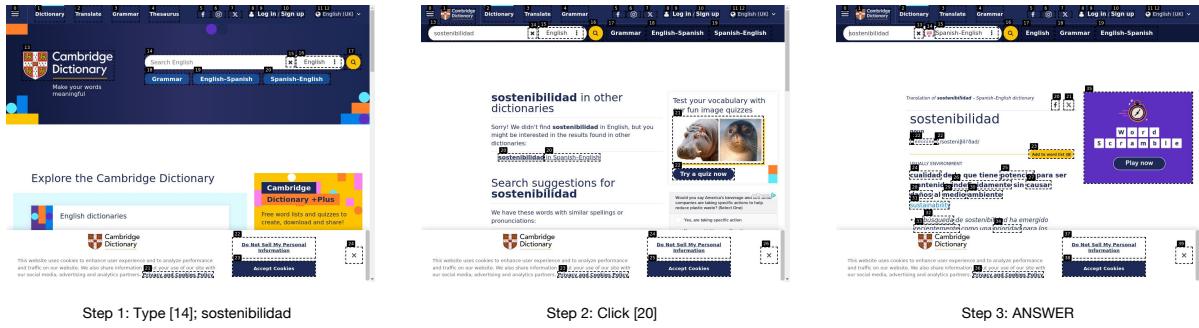
Figure 22: Screenshots of a complete trajectory of browsing Wolfram Alpha. Given the task: “Simplify  $x^5 - 20x^4 + 163x^3 - 676x^2 + 1424x - 1209$  so that it has fewer items.” The agent interacts with the Wolfram Alpha website and obtains the answer: “ $(x - 4)^5 + 3(x - 4)^3 + 7$ ”



Figure 23: Screenshots of a complete trajectory of browsing Google Flights in Chinese. Given the task: “Find a flight from Hangzhou to Shenzhen on June 15, 2024, returning on June 20, departing at night, answer a flight that meets the requirements, and answer its specific departure time, airline and price.” The agent interacts with the Google Flights website and obtains the answer: “Departure time: 17:35, airline: Shenzhen Airlines, price: HK\$2,680”

**TASK:** 请帮我查找一个航班，2024年6月15日杭州飞往深圳，6月20日返回，晚上出发，回答一个符合要求的航班，并回答其具体出发时间，航空公司以及价格。

**ANSWER:** 出发时间: 17:35, 航空公司: 深航, 价格: HK\$2,680



Step 1: Type [14]; sostenibilidad

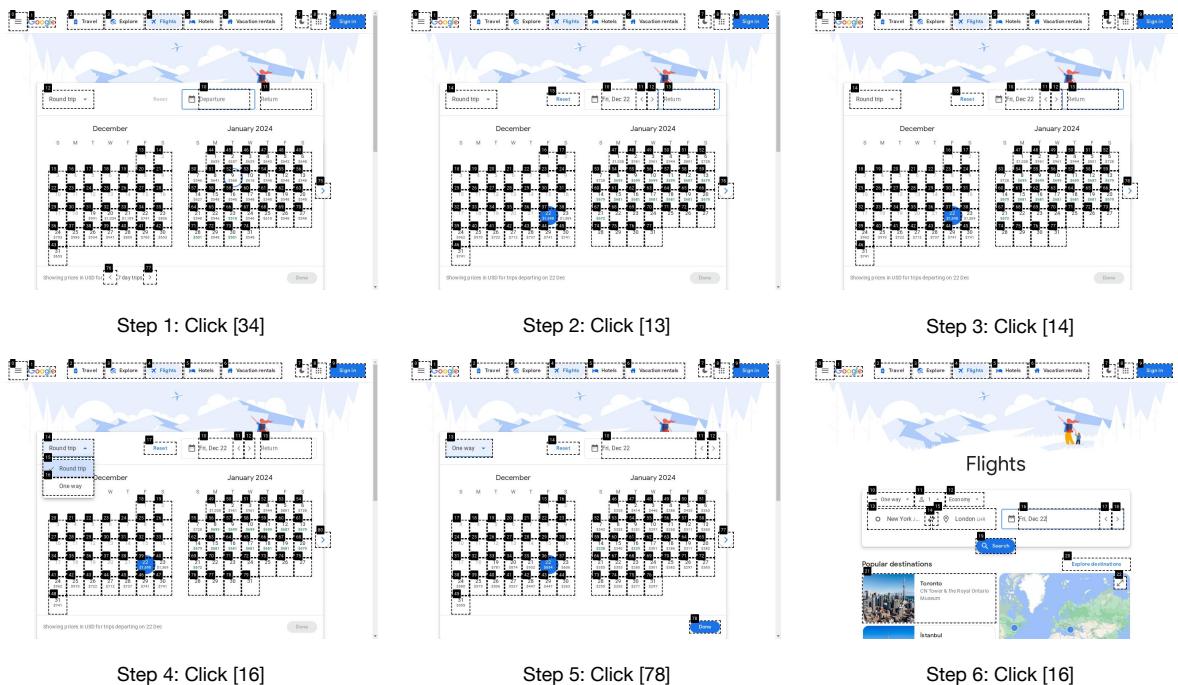
Step 2: Click [20]

Step 3: ANSWER

**TASK:** Busca la pronunciación y la definición de la palabra "sostenibilidad" en el diccionario.  
(Look up the pronunciation and definition of the word "sustainability" in the dictionary.)

**ANSWER:** The pronunciation of "sostenibilidad" is /sosteniβiliðað/ and it is a feminine noun. The definition relates to the quality of being able to continue over a period of time, causing minimal damage to the environment, which implies the term is equivalent to "sustainability" in English.

Figure 24: Screenshots of a complete trajectory of browsing Cambridge Dictionary in Spanish. The description of task and answer are shown in Figure.



Step 1: Click [34]

Step 2: Click [13]

Step 3: Click [14]

Step 4: Click [16]

Step 5: Click [78]

Step 6: Click [16]

Figure 25: An error case for Google Flights. Given the task: "Find the lowest fare from all eligible one-way flights for 1 adult from JFK to Heathrow on Jan. 22." Agent fails to select the correct numerical label though it really wants to select 22 January.

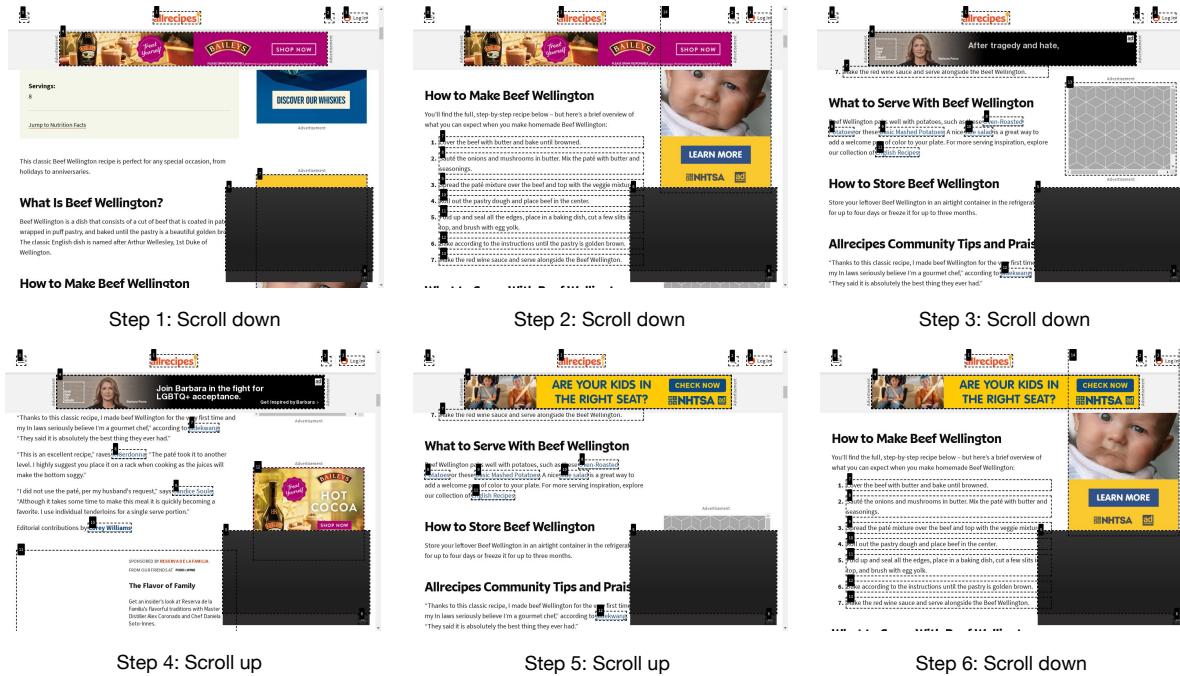


Figure 26: An error case for Allrecipes. Given the task: “Search for a recipe for Beef Wellington on Allrecipes that has at least 200 reviews and an average rating of 4.5 stars or higher. List the main ingredients required for the dish.” Agent fails to scroll the page correctly and find ingredients.

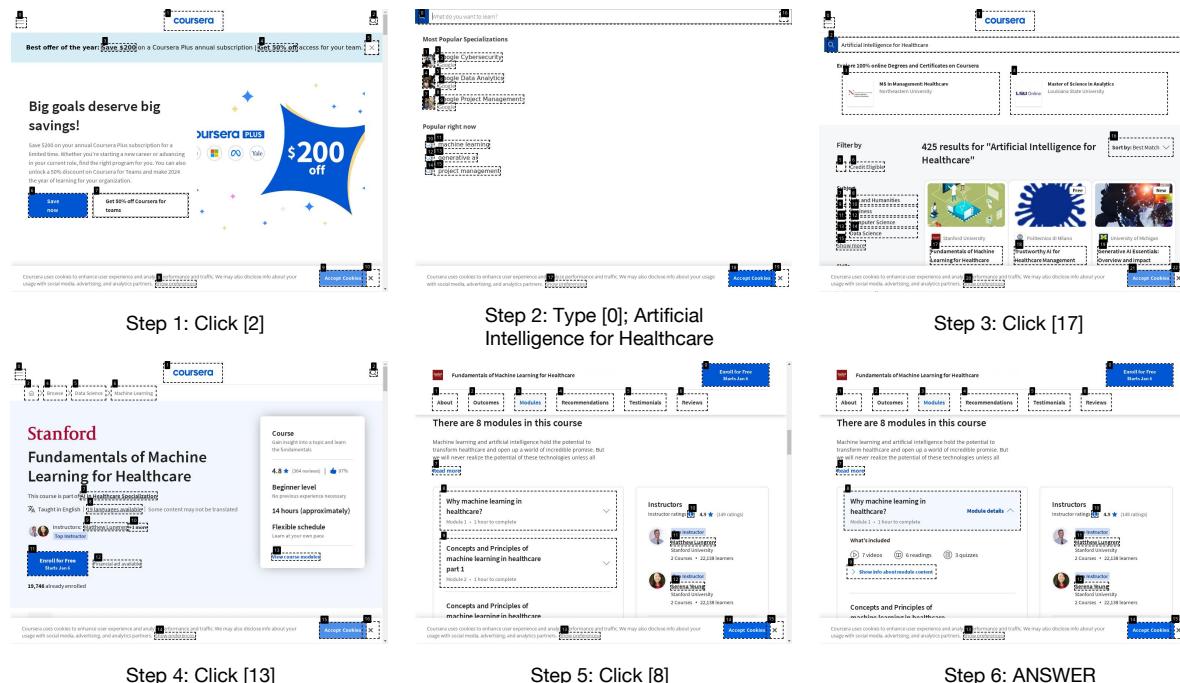


Figure 27: An error case for Coursera. Given the task: “Identify a course on Coursera related to ‘Artificial Intelligence for Healthcare’ and note the course duration along with the number of quizzes in Assessments.” The agent provides the answer: “Module 1 includes 3 quizzes in the Assessments.” However, we also require other Modules and agent fails to provide the results.



Figure 28: An error case for BBC News. Given the task: “Find out how many teams are in the Scottish Premiership of the Football Tournament and when did the Hibernian team’s most recent match start?” The agent provides the answer: “There are 12 teams in the Scottish Premiership. To find out the exact start time of Hibernian’s most recent match, further interaction with the website would be required.” The Agent knows that the task is not yet complete, but it ends its navigation early, even though it can find the Hibernian team’s most recent match.