
MIND2WEB: Towards a Generalist Agent for the Web

Xiang Deng* Yu Gu Boyuan Zheng Shijie Chen
Samuel Stevens Boshi Wang Huan Sun* Yu Su*
The Ohio State University

<https://osu-nlp-group.github.io/Mind2Web>

Abstract

We introduce MIND2WEB, the first dataset for developing and evaluating generalist agents for the web that can follow language instructions to complete complex tasks on any website. Existing datasets for web agents either use simulated websites or only cover a limited set of websites and tasks, thus not suitable for generalist web agents. With over 2,000 open-ended tasks collected from 137 websites spanning 31 domains and crowdsourced action sequences for the tasks, MIND2WEB provides three necessary ingredients for building generalist web agents: 1) diverse domains, websites, and tasks, 2) use of real-world websites instead of simulated and simplified ones, and 3) a broad spectrum of user interaction patterns. Based on MIND2WEB, we conduct an initial exploration of using large language models (LLMs) for building generalist web agents. While the raw HTML of real-world websites are often too large to be fed to LLMs, we show that first filtering it with a small LM significantly improves the effectiveness and efficiency of LLMs. Our solution demonstrates a decent level of performance, even on websites or entire domains the model has never seen before, but there is still a substantial room to improve towards truly generalizable agents. We open-source our dataset, model implementation, and trained models (<https://osu-nlp-group.github.io/Mind2Web>) to facilitate further research on building a generalist agent for the web.

1 Introduction

The web now hosts billions of websites [7] that cover virtually every aspect of the digital world. In this work, we seek to answer the question: *How can we build a generalist agent for the web that, given any website, can follow language instructions and carry out the corresponding tasks?* Some exemplar tasks for such an agent are shown in Figure 1. A generalist agent could make the web more accessible, which is becoming increasingly difficult as modern websites provide increasingly more functionalities that also increase their complexity and learning curve. On the other hand, such an agent may also turn the entire web into an unprecedentedly powerful and versatile tool [23, 27] that can enhance large language models (LLMs). For example, it may be used as a plugin for ChatGPT [24] to directly acquire information and carry out actions on HTML websites, instead of only retrieving web content through a retriever tool [8, 18] or relying on pre-defined APIs for each web service [34, 37].

A generalist agent for the web shall meet the following desiderata: First, **it shall work on any website on the Internet**. Since it is infeasible to collect sufficient training data that covers all websites, this requires the agent to be inherently generalizable to websites or even domains it has never seen before. Second, **it shall work on real-world websites, which can be dynamic, complex, and noisy**. Most modern websites are dynamic, generating and rendering different content in response to user actions. This necessitates the agent to model each website as a **partially-observable environment** instead of assuming full knowledge *a priori*. The agent should also not make strong simplifying assumptions

*Corresponding authors: {deng.595, sun.397, su.809}@osu.edu

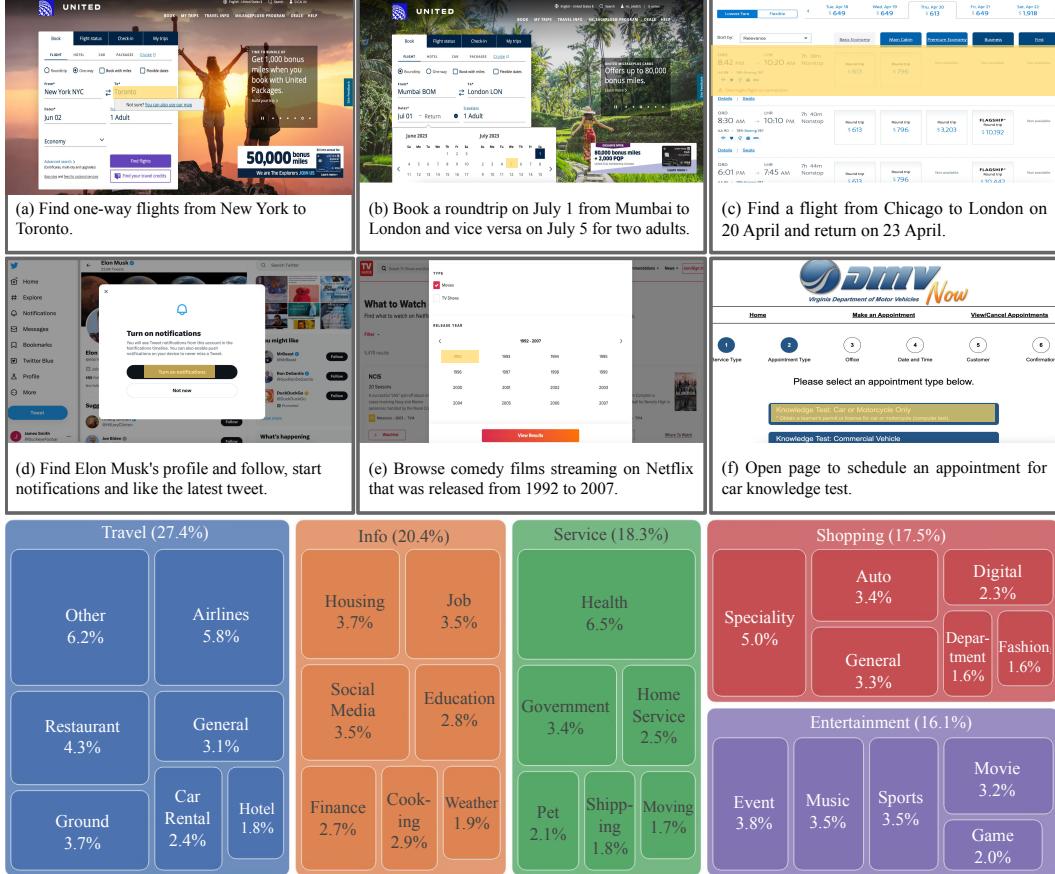


Figure 1: Sample tasks and all domains featured in MIND2WEB. The array of diversity allows for testing an agent’s generalizability across tasks on the same website (a vs. b), similar tasks on different websites (a vs. c), and even to entirely disparate tasks, websites, and domains (d–f).

about the environments but must embrace the full complexity and sometimes noise, *e.g.*, due to sub-optimal website designs. Finally, **it shall support diverse and sophisticated interactions with websites**. Tasks from users can be highly diverse and take a large number of steps to complete (*e.g.*, the task in Figure 1(b) would take 14 actions). An agent that only supports simple tasks may provide limited value to users.

Building an agent for the web is not entirely new. There have been numerous prior efforts in varied forms. However, none of them meets all the requirements for a generalist agent listed above. Existing work falls short in one to all of the following aspects: 1) Only operating in a limited and pre-specified set of websites [5 21 22 35 40], 2) making strong simplifying assumptions about the websites [22 40], and 3) only supporting specific types of tasks [21 22 40] and/or requiring tedious step-by-step instructions from users [5 21 22 39]. Meanwhile, LLMs have been shown to excel at grounded language understanding in complex environments with good generalizability and sample efficiency [2 13 17 33]. It is therefore promising to explore LLMs as a candidate solution towards generalist agents for the web. However, there lacks a good dataset that can support the development and evaluation of generalist web agents, which is the focus of this work.

In light of this, we present MIND2WEB, a new dataset with natural language tasks and manually annotated action sequences for developing and evaluating generalist agents for the web. It offers the following unique features:

1. Diverse coverage of domains, websites, and tasks. MIND2WEB boasts a collection of over 2,000 tasks curated from 137 websites that span 31 different domains. This extensive range of tasks and domains not only provides a vast landscape for exploration and learning but also opens up a new level of versatility and complexity, fostering a more comprehensive evaluation of generalist web agents.

2. Use of real-world websites. MIND2WEB replaces the oversimplified simulation environments commonly found in other datasets with an authentic, vibrant, and unpredictable realm of real-world websites. We provide full traces of user interactions, webpage snapshots, and network traffic, making it a rich source of raw, unfiltered, and dynamic data. By doing so, MIND2WEB equips models with the capacity to interact and cope with the complexities and uncertainties of real-world environments, thereby encouraging the development of more robust and adaptive models.

3. A broad spectrum of user interaction patterns. MIND2WEB enables users to engage in sophisticated ways with websites, as opposed to basic operations such as searching, following links, and reading content commonly found in existing work. Users can *click*, *select*, and *type in any elements* on the website, which significantly expands the space of possible tasks. This captures all the common actions users do on websites in real life and promotes the development of agents capable of handling complex tasks.

With the diverse domains and websites, we create challenging out-of-distribution evaluation settings where agents are tested on their generalizability to websites or even entire domains never seen during training. This presents a representative evaluation of generalist agents working on unseen websites.

MIND2WEB enables us to conduct an initial exploration in using LLMs for building generalist web agents. The HTML document of real-world webpages may contain thousands of elements, which makes it infeasible or too costly to be fed into an LLM’s context. To address this issue, we propose MINDACT, a two-stage model that involves first using a fine-tuned small LM to filter the web elements and then using an LLM to select from the filtered elements in a multi-choice question answering fashion and predict the corresponding action with the selected element. We show that MINDACT significantly outperforms modeling strategies commonly adopted by prior work and achieves a decent level of generalization. It can also work well with both open-source LLMs like Flan-T5 [31] through fine-tuning and closed-source LLMs like GPT-3.5-turbo and GPT-4 through in-context learning. However, there is still substantial room for further improvement towards generalist agents for the web. Promising future directions include integrating multi-modal information, reinforcement learning with feedback from real websites, and specialized LMs for web understanding and action taking.

2 MIND2WEB Dataset

Unlike existing datasets predominantly constructed within simulated environments [31] [40], our objective is to bridge the gap between simulation and reality so that agents trained on our dataset can work on real-world websites out of the box. To achieve this, our approach for data collection adheres to the following principles. Firstly, instead of recreating websites in simulation, which often leads to oversimplified environments, we engage directly with real-world websites and capture snapshots of these environments. Secondly, we collate a diverse set of websites from varied domains and crowd-source realistic tasks that cover a wide range of functionalities provided by these websites. Finally, acknowledging the challenge of perfectly replicating the complexity of real-world environments, we strive to capture a comprehensive snapshot of each website and the full interaction trace, to the extent that all the tasks can be seamlessly replayed offline. This supports rich modeling and evaluation approaches, ensuring a robust and practical dataset for research.

2.1 Task Definition

The primary objective of MIND2WEB is for the agent to complete a specific task on the target website through a series of actions. Each instance in our dataset contains three components:

Task description, which outlines the high-level goal of the task. We intentionally avoid low-level, step-by-step instructions, aiming to foster the development of agents that can comprehend and carry out tasks in a more autonomous fashion, rather than merely following prescriptive directives.

Action sequence, which is the sequence of actions required to accomplish the task on the website. Each action in the sequence comprises a (Target Element, Operation) pair. The Target Element is an interactable element on the current webpage, and the Operation refers to the action to be executed on that element. We support three common operations: Click (also including Hover and Press Enter), Type, and Select Option. For Type and Select Option, they also require an additional value as an argument. Actions in a sequence often span multiple webpages of a website.

Task Description:		Webpage Snapshots:		
<i>Show me the reviews for the auto repair business closest to 10002.</i>		Action 1	Action 2	Action 5
Action Sequence:				
Target Element	Operation	Action 1	Action 2	Action 5
1. [searchbox] <i>Find</i>	TYPE: <i>auto repair</i>			
2. [button] <i>Auto Repair</i>	CLICK			
3. [textbox] <i>Near</i>	TYPE: <i>10002</i>			
4. [button] <i>10002</i>	CLICK			
5. [button] <i>Search</i>	CLICK			
6. [switch] <i>Show BBB Accredited only</i>	CLICK			
7. [svg]	CLICK			
8. [button] <i>Sort By</i>	CLICK			
9. [link] <i>Fast Lane 24 Hour Auto Repair</i>	CLICK			
10. [link] <i>Read Reviews</i>	CLICK			

Figure 2: A sample data instance of our dataset with the three components. Actions marked in red will result in a transition to a new webpage.

Webpage snapshots, which constitute the environment within which the task is performed. We provide the snapshots in a variety of formats to accommodate different modeling approaches: self-contained MHTML file that includes the raw HTML code of the webpage, DOM snapshot containing the DOM tree along with the layout and style information of the screenshot of the rendered webpage, HAR file that includes all the network traffic for replaying the interaction if needed, and trace file that comprises the complete interaction trace during the task annotation process.

The agent receives the task description in the beginning. At each step, it also receives the current webpage and the history of previous actions. The objective is to accurately predict the subsequent action, which encompasses the target element for interaction and the operation.

2.2 Data Collection

Our data collection process consists of four stages: website selection, task proposal, task demonstration, and task verification. Website selection and task verification are done by the authors. For task proposal and demonstration, we develop a sophisticated annotation tool using Playwright² and hire annotators through Amazon Mechanical Turk. Refer to Supplementary for annotation tool details.

Website Selection. We start with 5 top-level domains: Travel, Shopping, Service, Entertainment, and Information, which are subsequently broken down into 31 (secondary) domains. We select websites within each domain based on their popularity in the US, as ranked by similarweb.com. We manually select 3-5 representative websites per domain, resulting in a collection of 137 websites in total.

Task Proposal. We present the annotators with a target website, a concise description of the website, and a few sample tasks associated with it. The annotators are then asked to propose *open-ended and realistic tasks* based on three criteria: the tasks should be of diverse types, require multiple rounds of interaction, and describe the high-level goal instead of step-by-step instructions. To further stimulate creativity and boost diversity, we use ChatGPT to generate seed tasks by prompting it to test different functionalities of a website. We generate 50 seed tasks per website, of which 10 are randomly sampled and presented to the annotator each time. These seed tasks are mainly for inspiration—annotators are explicitly instructed not to directly use them and we reject task proposals that are highly similar to the seed tasks. All the proposed tasks are further screened by the authors to ensure quality and diversity before entering the demonstration phase.

Task Demonstration. We develop a Playwright-based tool for demonstration (Figure 2). Workers will use the tool to demonstrate how to perform the tasks they have proposed within a web browser. To ensure accuracy, each interaction round is split into two parts: *element selection* and *operation selection*. At each step, the worker first selects an element on the webpage by clicking within the browser. They are then asked to confirm the selection and choose the operation to execute on the selected element. Once the task is completed, the worker is given another opportunity to review and modify the task description.

²<https://playwright.dev/>

Table 1: Statistics of MIND2WEB compared with existing datasets.

	# Dom.	# Env.	Env. Type	Avg. # Elements	# Tasks	Task Info.	Avg. # Actions
MiniWoB++ [22]	—	100	Simplified mobile websites	28	100	Low-level	3.6
WebShop [40]	1	1	Simplified shopping websites	38	12,000 products	High-level	11.3
RUSS [39]	—	22	Real-world websites	801	80	High & low	5.4
PixelHelp [21]	4	4	Mobile apps	—	187	High & low	-
META-GUI [35]	6	11	Mobile apps	79	1,125 dialogues	High-level	4.3
MoTIF [5]	15	125	Mobile apps	188	756	High & Low	4.4
MIND2WEB	5 / 31	137	Real-world websites	1,135	2,350	High-level	7.3

Task Verification. Lastly, all task demonstrations are verified by the authors to ensure the following: First, all actions are accurately reflected in the task description. The authors will modify the task description if needed to align it with the annotated actions; Second, the recorded actions are correct and clean, with extraneous steps discarded. Finally, the starting and ending points of tasks are consistent, such as excluding actions for closing popup windoes, or ending the annotation at the search result page if the task was to find a certain item without clicking on specific items. After verification, we discarded 61 out of the total 2,411 tasks. Among the 2,350 retained tasks, the task description was refined in 390 instances to better correspond with the demonstrated actions, while some extraneous steps were discarded in 187 instances. Overall, the data collection pipeline has been proven effective and produces high-quality data.

2.3 Comparison with Existing Work and Research Challenges

MIND2WEB presents a unique ensemble of research challenges for the development of generalist agents for the web in real-world settings. As shown in Table 1, MIND2WEB distinguishes itself from existing literature in several ways. Firstly, MIND2WEB spans across 137 websites from 31 domains, allowing comprehensive testing of an agent’s ability in generalizing across varied environments. Secondly, we utilize real-world websites without manual simplification. Consequently, the included environments exhibit complexity far surpassing that encountered in previous studies, yet better reflecting the intricacy of the modern web. With an average of over 1,000 elements per page embedded within complex DOM structures, how to effectively process such long and highly structured documents presents a significant challenge for modeling. Lastly, we direct the annotators to propose *open-ended tasks* that explore different functionalities of the website to mimic genuine web usage. Meanwhile, contrary to prior studies [5, 21, 22, 39] that provide step-by-step directives and primarily focus on testing the agent’s ability to translate low-level instructions into actions, *e.g.*, “*Type New York in the location field, click the search button and choose the tomorrow tab,*” we opted for the setting where only high-level goals are available, *e.g.*, “*What is the weather for New York tomorrow?*” This poses a much greater yet realistic planning and grounding challenge for the agent.

3 Method: MINDACT

Employing the data from MIND2WEB, we introduce an exploratory framework, MINDACT, for our task, leveraging the power of LLMs. Raw HTML documents, which could consist of thousands of elements, are either infeasible or cost-prohibitive to be directly fed into LLMs. We propose a two-stage process that synergizes the strength of small and large LMs, as shown in Figure 3. In the first stage, a fine-tuned small LM is used to rank the elements present on a webpage, yielding a small pool of promising candidates. In the second stage, these candidate elements are consolidated to form a representative snippet of the webpage, which is then processed by an LLM to predict the final action, including predicting both the element for interaction and the corresponding operation.

3.1 Candidate Generation with Small LMs

Given the task description, the snapshot of the webpage at step t , and the actions performed in the preceding $t - 1$ steps, we treat candidate generation as a ranking task. The task is to select the top- k

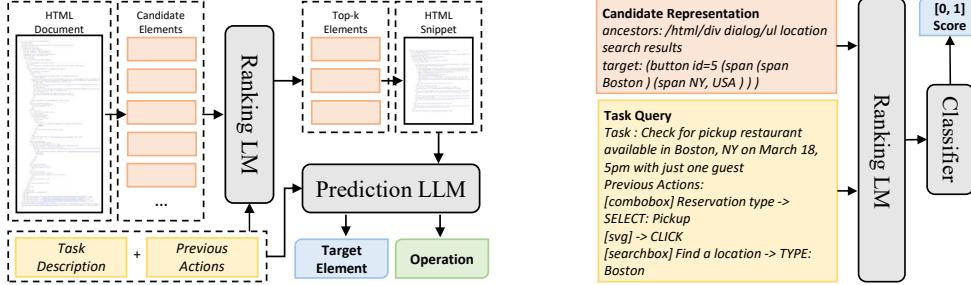


Figure 3: The overall pipeline for MINDACT with a small ranking LM for candidate generation, and a large prediction LM for action prediction.

Figure 4: Illustration of the candidate generation module and the templates for constructing task query and candidate representation.

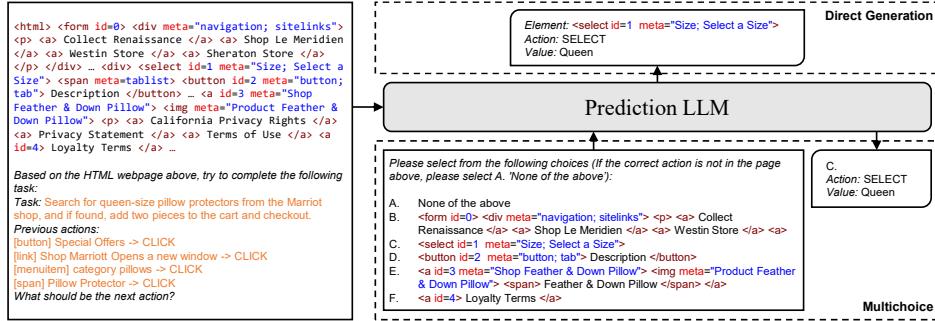


Figure 5: Illustration of action prediction with LLMs.

candidate DOM elements from the webpage that best align with both the task description and the current step. We formulate the task query by concatenating the task description with previous actions. The textual representation of each candidate DOM element is derived from a combination of the element’s tag, its textual content, and salient attribute values, as well as the textual representation of its parent and child elements. As shown in Figure 4, we pair each DOM element with the task query and feed it to an encoder-only LM through the cross-encoder architecture [28], yielding a matching score. At training time, we randomly sample negative elements from the webpage, and use the target element as the positive example. The matching score is passed through a sigmoid activation function and optimized with a binary cross entropy loss. At inference time, we score all elements in the webpage and pick the top- k elements with the largest logits as input to the second stage.

3.2 Action Prediction with LLMs

After obtaining the top- k candidates, we utilize the candidate set to prune the webpage snapshot and construct snippets that only include the selected candidates and their neighbours as inputs to an LLM. Recent studies [10, 13] have suggested that training LMs for discrimination rather than generation is more generalizable and sample-efficient for other grounding tasks. Inspired by that, we convert the task of element selection into a *multi-choice question answering* (QA) problem. Instead of generating the complete target element, the LM is trained to instead select from a list of options. For comparison, we also include a baseline that directly generates the target element based on the provided webpage snippet. In both cases, we directly let the LLM generate the operation, along with the additional value needed for some operations. An example is shown in Figure 5. We incorporate up to 5 candidate elements within each input, together with a None option, and partition the candidate set into several groups. During training, we construct the target sequence using ground-truth actions and fine-tune the model using a left-to-right language modeling objective. During inference, we divide the top- k candidates into multiple clusters of five options. If more than one option is selected after a round, we form new groups with the selected ones. This process repeats until a single element is selected, or all options are rejected by the model, i.e., the model chooses the None option for all groups.

Table 2: Main results. The classification baseline uses DeBERTa_B and the generation baseline uses Flan-T5_B. For step-wise metrics, we report macro average across tasks. * For GPT-4 we use 50 tasks for each setting with top-10 candidates due to limited budget. See Appendix D.3 for results on the 50 tasks subsets for all methods.

	Cross-Task				Cross-Website				Cross-Domain			
	Ele. Acc	Op. F1	Step SR	SR	Ele. Acc	Op. F1	Step SR	SR	Ele. Acc	Op. F1	Step SR	SR
Classification	26.8	—	—	—	21.6	—	—	—	24.5	—	—	—
Generation	20.2	52.0	17.5	0.0	13.9	44.7	11.0	0.0	14.2	44.7	11.9	0.4
MINDACT												
w/ Flan-T5 _B	43.6	76.8	41.0	4.0	32.1	67.6	29.5	1.7	33.9	67.3	31.6	1.6
w/ Flan-T5 _L	53.4	75.7	50.3	7.1	39.2	67.1	35.3	1.1	39.7	67.2	37.3	2.7
w/ Flan-T5 _{XL}	55.1	75.7	52.0	5.2	42.0	65.2	38.9	5.1	42.1	66.5	39.6	2.9
w/ GPT-3.5	20.3	56.6	17.4	0.8	19.3	48.8	16.2	0.6	21.6	52.8	18.6	1.0
w/ GPT-4*	41.6	60.6	36.2	2.0	35.8	51.1	30.1	2.0	37.1	46.5	26.4	2.0

4 Experiments

4.1 Experimental Setup

The diversity of MIND2WEB provides a unique opportunity to evaluate an agent’s generalizability at different levels. We seek to understand how well an agent can generalize across domains, websites, and tasks: **Test_{Cross-Domain}**, for which we hold out two top-level domains, *Information* and *Service*, with 912 tasks from 73 websites. Here, the model is expected to generalize to an entirely new domain without having seen any websites or tasks associated with that domain during training. **Test_{Cross-Website}**, with 10 websites from each remaining top-level domain, containing 177 tasks. In this setting, the model is never exposed to the test websites during training. However, it has been trained on websites from the same domain and possibly with similar tasks. This setup allows us to assess an agent’s ability to adapt to entirely new websites, yet within familiar domains and task contexts. **Test_{Cross-Task}**, where we randomly split 20% of the remaining data, regardless of domains and websites, resulting in 252 tasks from 69 websites. In this setting, the model has been exposed to webpages from the same website during training and has likely encountered similar tasks. The rest of data is used for training, which contains 1,009 tasks from 73 websites.

4.2 Data Preprocessing and Evaluation

We apply simple heuristics to clean the raw HTML documents, keeping only elements that are visible and carry substantial semantic meaning, as determined by their attributes, textual content, and neighboring elements. This effectively reduces the average number of elements from 1,135 to 580, while still maintaining an overall recall of 94.7% for the target element in the training data.

For evaluation, we first calculate **Element Accuracy** that compares the selected element with all acceptable elements, and **Operation F1** that calculates token-level F1 score for the predicted operation. This is the same as accuracy for Click, but considers the correctness of the input value for Type and Select Option. Each step of the task is evaluated independently with the ground truth action history provided. We then define **Step Success Rate** and **Success Rate** (for the whole task). A step is regarded as successful only if both the selected element and the predicted operation are correct. A task is regarded successful only if all steps have succeeded. It is therefore a stringent metric. For step-wise metrics, we report macro average across tasks.

4.3 Results

Candidate Generation. We fine-tune DeBERTa [16] as the small LM for candidate generation. As candidate generation requires high efficiency, we use the base version DeBERTa_B with 86M parameters. Overall, it achieves 88.9% / 85.3% / 85.7% Recall@50 on Test_{Cross-Task}, Test_{Cross-Website} and Test_{Cross-Domain}, respectively. We use its top-50 ranking results as the candidate pool for all subsequent experiments.

Action Prediction. We mainly compare against two baselines in Table 2. The first directly uses the candidate generation model (DeBERTa) for element selection, which is similar to existing work [14, 35] that combines an encoder with classification heads. However, such a design cannot

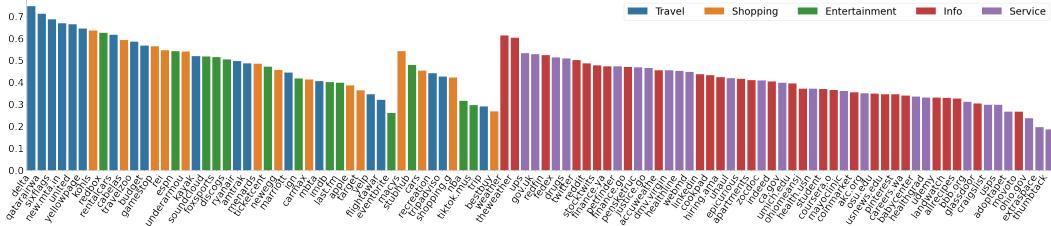


Figure 6: Step success rate per website grouped by the three splits: $\text{Test}_{\text{Cross-Task}}$, $\text{Test}_{\text{Cross-Website}}$ and $\text{Test}_{\text{Cross-Domain}}$ from left to right. Here we only show websites with more than three test tasks.

benefit from many recent LMs that use an encoder-decoder or decoder-only architecture. It cannot predict actions and the element selection performance is also not competitive, as shown in Table 2. We use Flan-T5 [10] as the backbone for the generation model. The autoregressive generation formulation (Figure 5 top) does not perform well, and even underperforms the classification baseline on element selection despite the larger model size (220M for Flan-T5_B). We observe a substantial gain with MINDACT using the multi-choice QA formulation. The best model achieves 52.0% step success rate under Cross-Task setting, and 38.9% / 39.6% when generalizing to unseen websites and domains. However, the overall task success rate remains low for all models, as the agent often commits at least one error step in most cases.

Three Levels of Generalization. All models perform best on the Cross-Task setting, with over 10% absolute gap (step SR) on average compared with Cross-Website and Cross-Domain settings, indicating that generalizing to unseen environments is still a major challenge. On the contrary, we note that the performance of Cross-Website and Cross-Domain settings are notably similar, which is also reinforced in Figure 6, where there is no clear distinction in performance across these settings. This suggests that the challenges primarily stem from the diversity in website designs and interaction logic rather than domain specifics. Tasks across domains tend to share common operations, and pretrained LMs may already have the capability to decompose complex tasks at a high level based on commonsense knowledge. Yet, grounding such knowledge into actionable steps in specific and varying environments remains a considerable challenge.

In-context Learning with LLMs. We also experiment with two popular LLMs, GPT-3.5-turbo and GPT-4 [25], through in-context learning. We use the same multiple-choice formulation as MINDACT, and include three demonstration examples for in-context learning. We can see that both models are comparable to the two baselines with only three in-context examples. Note that this is not a fair comparison with the Flan-T5 models, which are fine-tuned on the full training data. We also include the zero-shot results with Flan-T5_{XL} in Appendix D.2 but the model fails to perform the task without fine-tuning. Meanwhile, GPT-3.5 only has around 20% element selection accuracy, despite the superior performance people have observed on other datasets. Further analysis reveals that one possible problem is the model’s propensity to select the None option, asserting that the task cannot be finished on the current webpage. This is somewhat accurate since tasks typically necessitate navigation through multiple webpages and performing a series of actions before reaching the final result. This aspect indeed represents the primary difficulty of our task. On the other hand, we observe highly promising outcomes with GPT-4. The performance is on par with the tuned Flan-T5 models under Cross-Website and Cross-Domain settings for element selection, indicating a great potential for developing generalist agents using LLMs. Nevertheless, GPT-4’s high operational cost remains a concern. Developing smaller models specialized for the web is an interesting future avenue.

5 Related Work

Autonomous Agents for Web and Mobile Applications. Considerable initiatives have been invested in automating web navigation, driven by a vision of facilitating effortless human-web interaction. Yet, previous research has been limited by the types of tasks and websites it can handle, either confined to simplified simulation environments [22, 31, 40], or limited to a narrow set of domains and websites [39, 40]. Recent studies [5, 21, 35] have utilized similar techniques for mobile applications, however, these are often simpler and offer fewer functions compared with full-fledged websites. In contrast, MIND2WEB aims to adapt to a realistic web environment, characterized by its high diversity.

Also related is the research on web automation systems [1, 19]. These technologies often demand programming skills, which can make them less accessible to general users. We aim to equip the web automation system with a natural language interface, thereby reducing the entry barrier significantly.

Large Language Models. In recent years, there has been a surge in the development and application of large language models (LLMs). These models, often encompassing billions of parameters, are pre-trained on massive corpora of text data [3, 44, 45], enabling them to capture intricate linguistic patterns, nuances, and relationships, resulting in unprecedented performance on a wide array of NLP tasks. One of the most noteworthy attributes of LLMs is their few-shot learning capability. Unlike traditional machine learning models that necessitate extensive labeled data for task-specific fine-tuning, LLMs can often perform tasks with minimal task-specific examples. Furthermore, LLMs such as GPT-3 [4] and PaLM [9] have also demonstrated the ability to do in-context learning, where they can adapt to novel tasks by simply providing context within the input prompt, eliminating the need for explicit retraining. In this work, we explore the use of LLMs to build generalist agent on top of MIND2WEB by either tuning medium-sized LMs with only around 1,000 examples, or prompting an LLM such as GPT-4, and have observed promising results.

Grounded Language Understanding. Our work also aligns with the field of grounded language understanding, which aims to map natural language utterances onto executable plans in a target environment [13]. Many studies have centered around environments underpinned by a well-structured schema or ontology, including relational databases [36, 43] and knowledge bases [12, 42], which may not adequately reflect the more heterogeneous conditions in real-world situations. Our work instead grounds natural language in the noisy and schemaless web environment. Our setting is also connected to embodied AI, where an agent, guided by language instructions, carries out tasks in a physical environment [2, 32, 33]. Nonetheless, existing research primarily focuses on a specific setting (e.g., household environments), limiting their diversity. MIND2WEB provides a unique testbed for studying a broad range of grounding challenges in real-world environments.

Tool Learning. Recent developments have underscored LLMs’ potential in using a myriad of tools (*i.e.*, taking actions) to augment their capacity [23, 27], including search engine, translator, calculator, etc. Example works include Toolformer [29], ReAct [41], and ToolkenGPT [15]. The creation of recent benchmarks on tool learning [20, 26] further highlights the growing interest in evaluating LLMs’ proficiency in tool usage. However, existing research primarily concentrates on short-term tool invocation, neglecting long-term planning. MIND2WEB can bridge this lacuna by necessitating LLMs to take actions within realistic web-browsing environments that demand prolonged decision-making sequences. Furthermore, MIND2WEB may stimulate the development of more advanced tools based on LLMs that interface the web with natural language. These advanced tools could be subsequently employed by another LLM for more challenging problem-solving tasks [11, 30].

6 Limitations and Potential Societal Impact

MIND2WEB is designed to facilitate the development and evaluation of generalist agents for the web. Such agents hold great potential for making the web more accessible and easy to use, especially for individuals who are less familiar with information technology or have disabilities and may struggle to navigate through complex web apps and get overwhelmed by the options available. However, there are still potential concerns and limitations regarding the current data collection, system design and safety for deployment in real world.

Diversity and Representation in Data Collection. Although we strive to choose representative websites covering diverse domains, the present selection predominantly comprises English-language websites primarily used in the U.S. Meanwhile, all our annotators are sourced through the Amazon MTurk platform, which might be biased towards a group that is more proficient in web use. Therefore, the tasks and websites embodied in our dataset may represent only a subset of all potential tasks that can be performed on the web. Bearing this limitation in mind, the design of MIND2WEB and our data collection protocol allow for easy expansion to encompass more tasks and websites. The inclusion of additional websites, potentially from different countries and languages, and tasks from more diverse demographics, such as individuals from different age groups, those traditionally facing web accessibility challenges, and professionals from specific domains like software development, research, law, and more, present exciting directions for future development.

Use of Multimodal Information. Our current approach, MINDACT, models the web environment using only textual context from webpage snapshots. Nevertheless, crucial information can also be gleaned from the visual representation of a rendered webpage. While not currently utilized, we have included complete webpage snapshots in MIND2WEB, enabling rendering of the webpage for visual interpretation. The use of this multimodal information will be a viable prospect for improving model performance.

Modeling of Interaction Dynamics. In MINDACT, we encode each webpage independently at every step, with only the previous actions provided as historical context. However, the changes of the web environment could also provide significant cues for task completion, such as the appearance of a dropdown menu following a button click. Exploring effective ways to model such dynamic environment transformations during interaction could be an essential aspect for developing robust web agents.

Human-Agent Interaction. In the current design of MIND2WEB, the user provides a single description of the task goal up front, and the agent carries out the task from start to finish. In real-world settings, the user may wish to adjust or add task requirements in the middle, or the agent might seek user confirmation for more accurate task understanding. Extending Mind2Web to an interactive or conversational setting, thereby allowing diverse forms of human-agent interactions, could be an interesting future direction.

Evaluation with Offline/Online Environments. Following recent works [5, 35], we evaluate the system with cached offline environments, which allows us to test using snapshots of complex real-world websites. However, a downside to this is that the task will fail immediately if an action was not cached during data collection, potentially leading to false negatives due to the existence of multiple paths for completing the same task. As described in Appendix C.1, we normalize the actions to address equivalent elements within the same page. In addition, we include complete network traffic in the dataset, presenting possibilities for future research to enable some degree of replay and exploration within the cached environment. Given that MIND2WEB faithfully replicates real-world webpages, systems trained on the dataset should be readily transferable to live websites. Conducting end-to-end live evaluation on real websites with human assistance is a very promising direction that is worth exploration.

Safety in Deployment. While the development of general-purpose web agents holds great potential to enhance efficiency, optimize user experiences, and promote web accessibility universally, the accompanying safety considerations for real-world deployment cannot be ignored. These include how to effectively manage sensitive actions like financial transactions, enhancing transparency and interpretability, and keeping users in control during task execution. Additionally, there is the risk of these agents possessing the capability to breach existing security measures such as CAPTCHA and being exploited for malicious activities, such as disseminating false information. Therefore, it is also important for cybersecurity research to consider these potential uses and develop preemptive protective measures.

7 Conclusion

In this work, we introduced MIND2WEB, the first dataset for developing and evaluating generalist agents for the web. We also proposed MINDACT, an agent that leverages the power of (large) language models for effectively tackling this task. Our work opens up a wide range of promising future directions, including integrating multi-modal information, reinforcement learning with feedback from real websites, and specialized LMs for web understanding and action taking. We hope that MIND2WEB will serve as a valuable platform for the research community to advance towards generalist agents for the web.

Acknowledgements

The authors would thank colleagues from the OSU NLP group for constructive feedback and all contributors from the Amazon Mechanical Turk platform who participated in our study and assisted in data collection. This research was sponsored in part by NSF OAC 2112606, NSF CAREER #1942980, ARL W911NF2220144 and Ohio Supercomputer Center [6]. The views and conclusions contained herein are those of the authors and should not be interpreted as representing the official

policies, either expressed or implied, of the U.S. government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notice herein.

References

- [1] Puppeteer headless chrome node.js api. <https://github.com/puppeteer/puppeteer>, 2021.
- [2] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alexander Herzog, Daniel Ho, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Eric Jang, Rosario Jauregui Ruano, Kyle Jeffrey, Sally Jesmonth, Nikhil J. Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Kuang-Huei Lee, Sergey Levine, Yao Lu, Linda Luu, Carolina Parada, Peter Pastor, Jornell Quiambao, Kanishka Rao, Jarek Rettinghouse, Diego Reyes, Pierre Sermanet, Nicolas Sievers, Clayton Tan, Alexander Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, and Mengyuan Yan. Do as I can, not as I say: Grounding language in robotic affordances. *CoRR*, abs/2204.01691, 2022. doi: 10.48550/arXiv.2204.01691. URL <https://doi.org/10.48550/arXiv.2204.01691>.
- [3] Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ B. Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri S. Chatterji, Annie S. Chen, Kathleen Creel, Jared Quincy Davis, Dorottya Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah D. Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark S. Krass, Ranjay Krishna, Rohith Kuditipudi, and et al. On the opportunities and risks of foundation models. *CoRR*, abs/2108.07258, 2021. URL <https://arxiv.org/abs/2108.07258>.
- [4] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020.
- [5] Andrea Burns, Deniz Arsan, Sanjna Agrawal, Ranjitha Kumar, Kate Saenko, and Bryan A. Plummer. A dataset for interactive vision-language navigation with unknown command feasibility. In *European Conference on Computer Vision*, 2022.
- [6] Ohio Supercomputer Center. Ohio supercomputer center, 1987. URL <http://osc.edu/ark:/19495/f5s1ph73>.
- [7] Radoslav Chakarov. How many websites are there? how many are active in 2023? <https://webtribunal.net/blog/how-many-websites/>. 2023.
- [8] Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. Reading Wikipedia to answer open-domain questions. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1870–1879, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi: 10.18653/v1/P17-1171. URL <https://aclanthology.org/P17-1171>.
- [9] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope,

- James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. Palm: Scaling language modeling with pathways, 2023. URL <http://jmlr.org/papers/v24/22-1144.html>.
- [10] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Y. Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models. *CoRR*, abs/2210.11416, 2022. doi: 10.48550/arXiv.2210.11416. URL <https://doi.org/10.48550/arXiv.2210.11416>.
- [11] Yingqiang Ge, Wenyue Hua, Jianchao Ji, Juntao Tan, Shuyuan Xu, and Yongfeng Zhang. Openagi: When LLM meets domain experts. *CoRR*, abs/2304.04370, 2023. doi: 10.48550/arXiv.2304.04370. URL <https://doi.org/10.48550/arXiv.2304.04370>.
- [12] Yu Gu, Sue Kase, Michelle Vanni, Brian M. Sadler, Percy Liang, Xifeng Yan, and Yu Su. Beyond I.I.D.: three levels of generalization for question answering on knowledge bases. In Jure Leskovec, Marko Grobelnik, Marc Najork, Jie Tang, and Leila Zia, editors, *WWW '21: The Web Conference 2021, Virtual Event / Ljubljana, Slovenia, April 19-23, 2021*, pages 3477–3488. ACM / IW3C2, 2021. doi: 10.1145/3442381.3449992. URL <https://doi.org/10.1145/3442381.3449992>.
- [13] Yu Gu, Xiang Deng, and Yu Su. Don't generate, discriminate: A proposal for grounding language models to real-world environments. *CoRR*, abs/2212.09736, 2022. doi: 10.48550/arXiv.2212.09736. URL <https://doi.org/10.48550/arXiv.2212.09736>.
- [14] Izzeddin Gur, Ofir Nachum, Yingjie Miao, Mustafa Safdari, Austin Huang, Aakanksha Chowdhery, Sharan Narang, Noah Fiedel, and Aleksandra Faust. Understanding html with large language models, 2023.
- [15] Shibo Hao, Tianyang Liu, Zhen Wang, and Zhiting Hu. Toolkengpt: Augmenting frozen language models with massive tools via tool embeddings. *CoRR*, abs/2305.11554, 2023. doi: 10.48550/arXiv.2305.11554. URL <https://doi.org/10.48550/arXiv.2305.11554>.
- [16] Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: decoding-enhanced bert with disentangled attention. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. URL <https://openreview.net/forum?id=XPZIaotutsD>.
- [17] Jinhao Jiang, Kun Zhou, Zican Dong, Keming Ye, Wayne Xin Zhao, and Ji-Rong Wen. Structgpt: A general framework for large language model to reason over structured data. *CoRR*, abs/2305.09645, 2023. doi: 10.48550/arXiv.2305.09645. URL <https://doi.org/10.48550/arXiv.2305.09645>.
- [18] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.550. URL <https://aclanthology.org/2020.emnlp-main.550>.
- [19] Gilly Leshed, Eben M. Haber, Tara Matthews, and Tessa Lau. CoScripter: Automating & sharing how-to knowledge in the enterprise. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 1719–1728, Florence Italy, April 2008. ACM. ISBN 978-1-60558-011-1. doi: 10.1145/1357054.1357323.

- [20] Minghao Li, Feifan Song, Bowen Yu, Haiyang Yu, Zhoujun Li, Fei Huang, and Yongbin Li. Api-bank: A benchmark for tool-augmented llms. *CoRR*, abs/2304.08244, 2023. doi: 10.48550/arXiv.2304.08244. URL <https://doi.org/10.48550/arXiv.2304.08244>.
- [21] Yang Li, Jiacong He, Xin Zhou, Yuan Zhang, and Jason Baldridge. Mapping natural language instructions to mobile UI action sequences. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel R. Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 8198–8210. Association for Computational Linguistics, 2020. doi: 10.18653/v1/2020.acl-main.729. URL <https://doi.org/10.18653/v1/2020.acl-main.729>.
- [22] Evan Zheran Liu, Kelvin Guu, Panupong Pasupat, Tianlin Shi, and Percy Liang. Reinforcement learning on web interfaces using workflow-guided exploration. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL <https://openreview.net/forum?id=ryTp3f-0->
- [23] Grégoire Mialon, Roberto Dessì, Maria Lomeli, Christoforos Nalmpantis, Ramakanth Pasunuru, Roberta Raileanu, Baptiste Rozière, Timo Schick, Jane Dwivedi-Yu, Asli Celikyilmaz, Edouard Grave, Yann LeCun, and Thomas Scialom. Augmented language models: a survey. *CoRR*, abs/2302.07842, 2023. doi: 10.48550/arXiv.2302.07842. URL <https://doi.org/10.48550/arXiv.2302.07842>.
- [24] OpenAI. Chatgpt plugins. <https://openai.com/blog/chatgpt-plugins>, 2023.
- [25] OpenAI. Gpt-4 technical report, 2023.
- [26] Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. Gorilla: Large language model connected with massive apis. *CoRR*, abs/2305.15334, 2023. doi: 10.48550/arXiv.2305.15334. URL <https://doi.org/10.48550/arXiv.2305.15334>.
- [27] Yujia Qin, Shengding Hu, Yankai Lin, Weize Chen, Ning Ding, Ganqu Cui, Zheni Zeng, Yufei Huang, Chaojun Xiao, Chi Han, Yi Ren Fung, Yusheng Su, Huadong Wang, Cheng Qian, Runchu Tian, Kunlun Zhu, Shihao Liang, Xingyu Shen, Bokai Xu, Zhen Zhang, Yining Ye, Bowen Li, Ziwei Tang, Jing Yi, Yuzhang Zhu, Zhenning Dai, Lan Yan, Xin Cong, Yaxi Lu, Weilin Zhao, Yuxiang Huang, Junxi Yan, Xu Han, Xian Sun, Dahai Li, Jason Phang, Cheng Yang, Tongshuang Wu, Heng Ji, Zhiyuan Liu, and Maosong Sun. Tool learning with foundation models. *CoRR*, abs/2304.08354, 2023. doi: 10.48550/arXiv.2304.08354. URL <https://doi.org/10.48550/arXiv.2304.08354>.
- [28] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019. URL <https://arxiv.org/abs/1908.10084>.
- [29] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *CoRR*, abs/2302.04761, 2023. doi: 10.48550/arXiv.2302.04761. URL <https://doi.org/10.48550/arXiv.2302.04761>.
- [30] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving AI tasks with chatgpt and its friends in huggingface. *CoRR*, abs/2303.17580, 2023. doi: 10.48550/arXiv.2303.17580. URL <https://doi.org/10.48550/arXiv.2303.17580>.
- [31] Tianlin Shi, Andrej Karpathy, Linxi Fan, Jonathan Hernandez, and Percy Liang. World of Bits: An Open-Domain Platform for Web-Based Agents. In *Proceedings of the 34th International Conference on Machine Learning*, pages 3135–3144. PMLR, July 2017.
- [32] Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. ALFRED: A benchmark for interpreting grounded instructions for everyday tasks. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 10737–10746.

- Computer Vision Foundation / IEEE, 2020. doi: 10.1109/CVPR42600.2020.01075. URL https://openaccess.thecvf.com/content_CVPR_2020/html/Shridhar_ALFRED_A_Benchmark_for_Interpreting_Grounded_Instructions_for_Everyday_Tasks_CVPR_2020_paper.html.
- [33] Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M. Sadler, Wei-Lun Chao, and Yu Su. Llm-planner: Few-shot grounded planning for embodied agents with large language models. *CoRR*, abs/2212.04088, 2022. doi: 10.48550/arXiv.2212.04088. URL <https://doi.org/10.48550/arXiv.2212.04088>.
- [34] Yu Su, Ahmed Hassan Awadallah, Madian Khabsa, Patrick Pantel, Michael Gamon, and Mark J. Encarnación. Building natural language interfaces to web apis. In Ee-Peng Lim, Marianne Winslett, Mark Sanderson, Ada Wai-Chee Fu, Jimeng Sun, J. Shane Culpepper, Eric Lo, Joyce C. Ho, Debora Donato, Rakesh Agrawal, Yu Zheng, Carlos Castillo, Aixin Sun, Vincent S. Tseng, and Chenliang Li, editors, *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management, CIKM 2017, Singapore, November 06 - 10, 2017*, pages 177–186. ACM, 2017. doi: 10.1145/3132847.3133009. URL <https://doi.org/10.1145/3132847.3133009>.
- [35] Liangtai Sun, Xingyu Chen, Lu Chen, Tianle Dai, Zichen Zhu, and Kai Yu. META-GUI: Towards Multi-modal Conversational Agents on Mobile GUI, November 2022.
- [36] Bailin Wang, Richard Shin, Xiaodong Liu, Oleksandr Polozov, and Matthew Richardson. RAT-SQL: Relation-Aware Schema Encoding and Linking for Text-to-SQL Parsers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7567–7578, Online, 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.677.
- [37] Kyle Williams, Seyyed Hadi Hashemi, and Imed Zitouni. Automatic task completion flows from web apis. In Benjamin Piwowarski, Max Chevalier, Éric Gaussier, Yoelle Maarek, Jian-Yun Nie, and Falk Scholer, editors, *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2019, Paris, France, July 21-25, 2019*, pages 1009–1012. ACM, 2019. doi: 10.1145/3331184.3331318. URL <https://doi.org/10.1145/3331184.3331318>.
- [38] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierrick Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online, October 2020. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/2020.emnlp-demos.6>.
- [39] Nancy Xu, Sam Masling, Michael Du, Giovanni Campagna, Larry Heck, James A. Landay, and Monica S. Lam. Grounding open-domain instructions to automate web support tasks. In *North American Chapter of the Association for Computational Linguistics*, 2021.
- [40] Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. WebShop: Towards Scalable Real-World Web Interaction with Grounded Language Agents. July 2022. doi: 10.48550/arXiv.2207.01206.
- [41] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. *CoRR*, abs/2210.03629, 2022. doi: 10.48550/arXiv.2210.03629. URL <https://doi.org/10.48550/arXiv.2210.03629>.
- [42] Wen-tau Yih, Ming-Wei Chang, Xiaodong He, and Jianfeng Gao. Semantic Parsing via Staged Query Graph Generation: Question Answering with Knowledge Base. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1321–1331, Beijing, China, 2015. Association for Computational Linguistics. doi: 10.3115/v1/P15-1128.

- [43] Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir R. Radev. Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task. In Ellen Riloff, David Chiang, Julia Hockenmaier, and Jun’ichi Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 3911–3921. Association for Computational Linguistics, 2018. doi: 10.18653/v1/d18-1425. URL <https://doi.org/10.18653/v1/d18-1425>.
- [44] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. A survey of large language models. *CoRR*, abs/2303.18223, 2023. doi: 10.48550/ARXIV.2303.18223. URL <https://doi.org/10.48550/arXiv.2303.18223>.
- [45] Ce Zhou, Qian Li, Chen Li, Jun Yu, Yixin Liu, Guangjing Wang, Kai Zhang, Cheng Ji, Qiben Yan, Lifang He, Hao Peng, Jianxin Li, Jia Wu, Ziwei Liu, Pengtao Xie, Caiming Xiong, Jian Pei, Philip S. Yu, and Lichao Sun. A comprehensive survey on pretrained foundation models: A history from BERT to chatgpt. *CoRR*, abs/2302.09419, 2023. doi: 10.48550/ARXIV.2302.09419. URL <https://doi.org/10.48550/arXiv.2302.09419>.

A Overview

Our supplementary includes the following sections:

- **Section B: Data Collection Details.** Details for crowdsourcing and implementation details for the three data collection phases: task proposal, task demonstration, and task verification.
- **Section C: Experiment Details.** Details for evaluation and model implementation.
- **Section D: Additional Results.** Results for additional auxiliary experiments.

Following NeurIPS Dataset and Benchmark track guidelines, we have shared the following artifacts:

Artifcat	Link	License
Homepage	https://osu-nlp-group.github.io/Mind2Web/	-
Code Repository	https://github.com/OSU-NLP-Group/Mind2Web	MIT License
Training Data	https://huggingface.co/datasets/osunlp/Mind2Web	CC BY 4.0
Test Data	https://shorturl.at/iGI45 (password “mind2web”)	CC BY 4.0

The authors and the OSU NLP group are committed to ensuring its regular upkeep and updates.

Table 3: Prompt for ChatGPT to generate seed tasks to inspire the annotators during task proposal.

Prompt	We are testing the website https://aa.com American Airlines - Airline tickets and low fares at aa.com Book low fares to destinations around the world and find the latest deals on airline tickets, hotels, car rentals and vacations at aa.com. As an AAdvantage member you earn miles on every trip and everyday spend. We want to test practical daily tasks that a user would do on the website. Come up with a list of 5 example tasks and try to cover different cases. Requirements: <ul style="list-style-type: none"> - Each example should be a single sentence and not just click one of the elements. - Don't give step-by-step instructions or directly mention the element to interact. - Describe the goal of the task and provide concrete information or constraints. Use mock-up information (identifier, number, personal information, name, date, attributes, constraints, etc.) to make the task more specific and realistic.
Sample Responses	Retrieve the confirmation number for a reservation made under the name jane smith departing from dallas (dfw) to miami (mia) on january 20th. Find and book a round-trip flight from new york to london departing on december 15th and returning on december 22nd for under \$800. Find the lowest fare for a one-way flight from new york city (jfk) to los angeles (lax) leaving on august 15th and returning on august 20th, with no layovers exceeding 2 hours. Rent a car in honolulu for one week starting on august 1st that can fit at least four passengers and has gps navigation included. Cancel a car rental reservation for a compact car in orlando, florida that was scheduled to be picked up on may 25th at 5 pm without incurring any cancellation fees.

B Data Collection Details

B.1 Crowdsourcing Details

In this study, we employ annotators from the Amazon Mechanical Turk (mturk) platform. The annotators are required to have a minimum of 1,000 approved HITs with an approval rate exceeding 98% on the platform. We design the compensation with an estimated hourly rate of \$10.10 to respect the minimum wage guidelines in Ohio, United States. Every worker passing our qualification receives a bonus, and we pay \$0.80 for each approved final task. We do not collect any identifiable private information during the study, and explicitly instruct the annotators to refrain from entering personal or sensitive data into the system. Annotators engage with our annotation tool only within a secure, remote sandbox environment, posing no foreseeable harm. The study complies with the IRB exemption criteria, per the Office of Responsible Research Practices at The Ohio State University. All annotators are presented with a consent form, to which they must agree before participating in the study. To prepare the workers for the task, we provide a comprehensive training document and a video tutorial, followed by a qualification assessment comprising a questionnaire and a series of test demonstrations using our tool. It is noteworthy that the task is divided into two phases: task proposal and task demonstration. The proposal phase comes with a nominal reward, with the majority of the compensation dispensed upon successful completion of the demonstration.

Quality and diversity are ensured through a two-stage review process. The first author reviews all tasks after task proposal and manually select the tasks for demonstration. After task demonstration, a thorough final verification of all collected data is conducted by all authors to authenticate the tasks and recorded actions. Each demonstration is first verified by one of the authors, and uncertain ones are further verified by the first author to reach consensus.

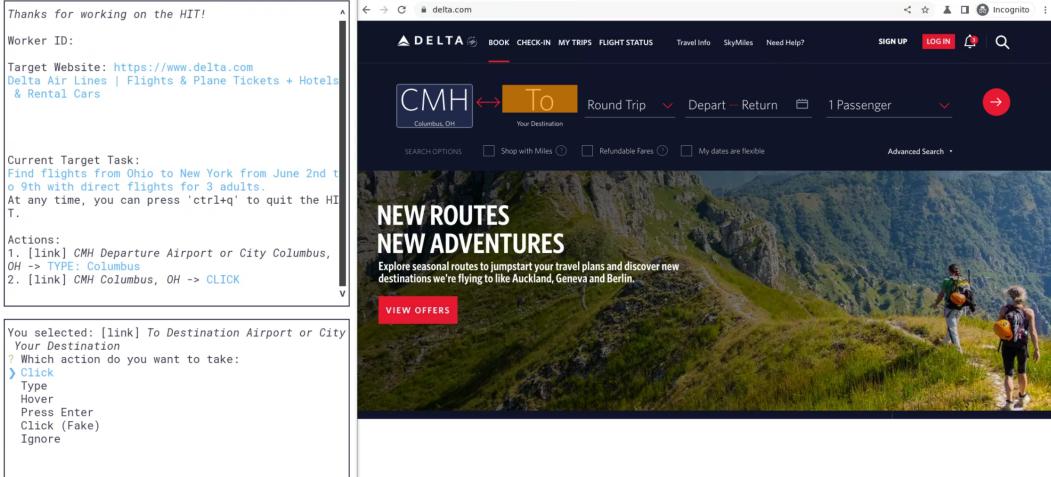


Figure 7: Illustration of our annotation tool, which consists of two side-by-side windows. On the left we provide a dialogue window for the user to control the tool and select operations to take. On the right we provide the browser window for the user to interact with and select web elements.

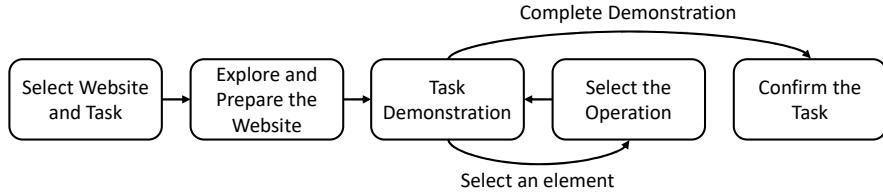


Figure 8: The overall procedure for task demonstration.

B.2 Task Proposal

We use ChatGPT to generate sample tasks to provide inspiration to the annotators, and the prompt used is shown in Table 3. For each HIT, we ask the annotator to select a website of their interest first. Following this, we present them with ten sample tasks produced by ChatGPT, and request them to propose a maximum of five additional tasks. The annotator is instructed not to directly copy the sample tasks. We manually evaluate all submitted tasks and reject those that demonstrate low quality or are too similar to previously accepted tasks. We set a nominal reward of \$0.05 for each task proposal HIT, and the annotator will receive it no matter whether the tasks are accepted or not. For accepted ones, the annotator will receive the full reward of \$0.80 on successful demonstration of the task. Once we have collected a total of around 20 tasks for a specific website, we desist from showing it to the user, aiming for a balanced distribution among websites and increased task diversity.

B.3 Task Demonstration

We develop a dedicated annotation tool for task demonstration using Playwright³ which allows us to interact with the browser and record user actions. As shown in Figure 7, the tool is composed of two windows. The dialogue window on the left serves as the annotator's control panel for guiding the interaction flow and choosing operations. The browser window on the right is where the annotator navigates the website and selects elements for interaction. Figure 8 shows the overall procedure for task demonstration. The annotator starts by selecting the website and task to be demonstrated. Once selected, the tool will bring up the website in the browser. The annotator is then instructed to explore the website and practice the task. To collect clean actions during actual demonstration, the workers are asked to close pop-up windows during exploration. We also provide anonymous accounts for the workers to use so that no private information is entered. The exploration stage is

³<https://playwright.dev/>

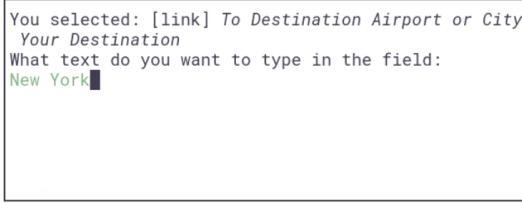


Figure 9: Dialogue window for the Type operation.

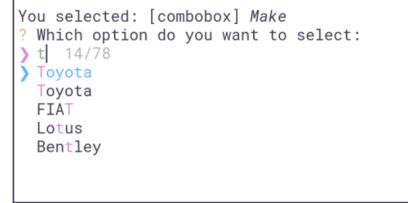


Figure 10: Dialogue window for the Select Option operation.

not recorded and primarily serves to familiarize the annotator with the website and task, as well as to prepare the website to prevent future pop-ups, thereby ensuring a clean, streamlined set of final recorded actions. After exploration, the annotator is directed to return to the homepage and reset any altered values, allowing us to begin the demonstration in a fresh state. During the demonstration, the annotator will illustrate how to accomplish the task step-by-step using both the browser and the dialogue window. To ensure a clean set of annotated actions, annotators are restricted from directly engaging with the browser during the demonstration phase. Instead, we divide each action step into two stages: Element selection and operation selection. At each step, the annotator first selects the target element by clicking it in the browser. We will highlight the selected element in the browser window but block the actual click event. The annotator is then prompted to select the operation to perform within the dialogue window, which is then carried out by the annotation tool in the browser. We provide 6 operations: Click, Type, Hover, Press Enter, Click (Fake) and Ignore. For the Type operation, the annotator is additionally required to supply the value as shown in Figure 9. If the chosen element is a select HTML element, and the annotator opts for Click, it translates to a Select Option operation and we will prompt the annotator to select one of the options as shown in Figure 10. To avoid ambiguity, the Click, Hover and Press Enter operations are all mapped to Click in the final dataset. Click (Fake) is a special operation. It will be recorded the same as a normal Click but will not get executed in the browser. This is designed for safeguarding against state-changing actions (i.e., actions that produce side effects to the world), such as posting a comment or scheduling an appointment, since it will interfere with other real users of the website. In practice, once a model predicts Click (Fake), it may prompt the user for confirmation before executing such state-changing actions. Finally, the annotator can also choose Ignore in case they select a wrong element. Once all the actions have been annotated, the annotator can choose to complete the task. They will then be asked to confirm the task description again and make any necessary modifications.

Pop-ups and CAPTCHAs. In this study, we emphasize on clean and direct task execution actions, intentionally omitting extraneous steps like pop-ups and CAPTCHAs that might introduce ambiguity in evaluation. We carefully select only those websites that pose no access issues when used with our tool. Before recording the task demonstration, annotators are requested to familiarize themselves with the website, and preemptively close pop-up windows and clear CAPTCHAs to avoid their recurrence during the actual demonstration. Annotators are further guided not to engage in extra steps such as closing ads unless necessary during the task demonstration. In the final task verification, we revisit the actions and filter out those unrelated to direct task execution. At the same time, we acknowledge that these instances constitute a significant aspect of the dynamic web environment in the real world. Enhancing systems to robustly tackle such scenarios on-the-go could form an interesting avenue for future research.

Mitigating Disruptions on the Websites. The annotator are advised against actions that could potentially interfere the normal operation of the website. To handle tasks such as scheduling appointments, we introduce a Click (Fake) operation that annotators can utilize to indicate the action without actually executing it on the website.

B.4 Task Verification

All collected data undergoes an additional verification process conducted by the authors, as demonstrated in Figure 11. The verification interface is shown in Figure 11. This verification consists of three tasks. Firstly, we evaluate whether a task should be discarded due to its low quality. Secondly, we examine each step to determine if any action should be discarded. This includes reviewing the

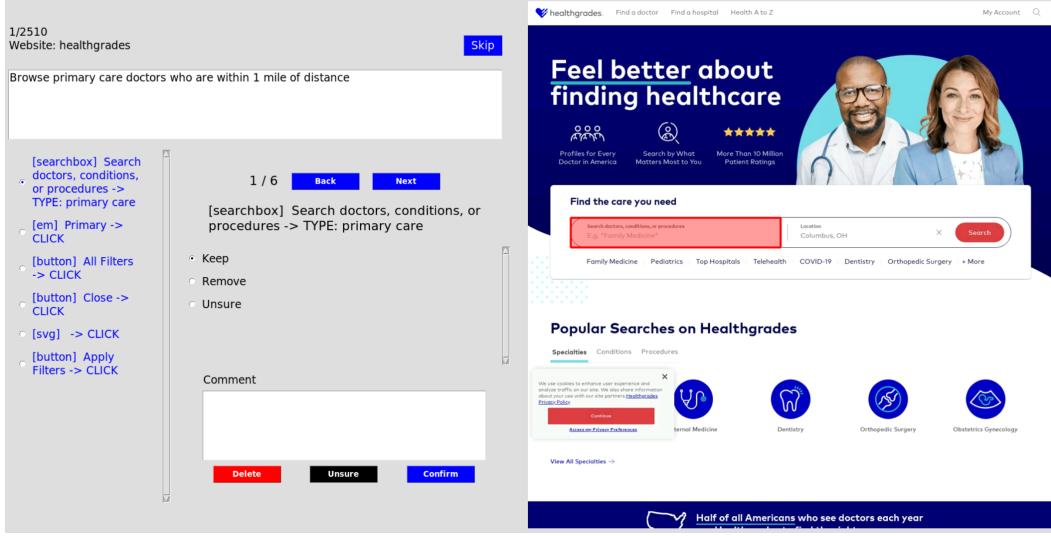


Figure 11: Illustration of our verification tool.

initial and final actions of the task, and excluding any additional actions (e.g., closing ads) that are not outlined in the task description to ensure consistency across task annotations. Finally, we verify the task description to confirm that all actions are accurately represented and make modifications if necessary. If there is uncertainty regarding any action, the verifier can opt for the ‘unsure’ option, prompting a re-evaluation by the first author.

C Experiment Details

C.1 Evaluation

One complication that arises during evaluation on real-world websites is that multiple elements on a webpage may induce the same effect. For instance, a button might house a text span within it, both of which, when clicked, yield identical results. To enhance the robustness of our evaluation, we employ heuristics to detect elements equivalent to the ground truth. We first examine the ancestors of the labeled element to identify potential higher-level elements acceptable for the current action. We employ a straightforward heuristic that locates the nearest clickable element to the ground truth, including itself. After identifying the top-level acceptable element, we include all its visible descendants that are located within its post-rendering bounding box as acceptable as well. Manual checking on 100 instances where the heuristic identifies a top-level element other than the ground truth confirms the validity of the approach. For both training and evaluation stages, all acceptable elements are considered positive.

C.2 Model Implementation Details

Candidate Generation. We use the Cross-Encoder implementation from Sentence-Transformers⁴ and use DeBERTa as the backbone model. More specifically, we use DeBERTa-v3-base⁵ for our experiments.

Action Prediction. We use the Seq2Seq model implementation from Transformers⁶. We experiment with the base⁶, large⁷ and xl⁸ versions of Flan-T5⁹.

⁴<https://www.sbert.net/examples/applications/cross-encoder/README.html>
⁵<https://huggingface.co/microsoft/deberta-v3-base>
⁶<https://huggingface.co/google/flan-t5-base>
⁷<https://huggingface.co/google/flan-t5-large>
⁸<https://huggingface.co/google/flan-t5-xl>

Table 4: Hyperparameters used in experiments.

Method	Model	Hyperparamerters
Candidate Generation	deberta-v3-base	batch_size:32, epoch:5, learning_rate:3e-5
Action Prediction Generation	flan-t5-base	batch_size:32, epoch:5, learning_rate:5e-5
MINDACT	flan-t5-base, flan-t5-large, flan-t5-xl gpt-3.5-turbo, gpt-4	batch_size:32, epoch:5, learning_rate:5e-5 temperature:0, # demonstrations:3

Table 5: Step Success Rate for Flan-T5 models with different groups of options. Here we show mean and standard deviation of 5 runs with different random seeds.

	Cross-Task	Cross-Website	Cross-Domain
Flan-T5 _B	41.5±0.7	30.0±0.8	31.3±0.5
Flan-T5 _L	49.9±0.2	35.7±0.5	36.7±0.3
Flan-T5 _{XL}	51.9±0.8	39.5±0.2	39.6±0.2

LLM In-context Learning. We use the OpenAI API for in-context learning with LLMs. We experiment with two versions of GPT models: gpt-3.5-turbo and gpt-4. We include three demonstration examples for in-context learning. The complete prompt is shown in Table 8.

Training Details. The flan-t5-xl and flan-t5-large models are trained on servers with 4*A100 80GB cards provided by Ohio Supercomputer Center [6]. All other models are trained with single A6000 48GB cards.

Please see Table 4 for all hyperparameters used in our experiments.

D Additional Results

D.1 Effect of Random Grouping Elements for Action Prediction

For both training and inference, we shuffle the elements in the webpage and randomly group them into multi-choice questions. The model might give different predictions when presented with different sets of choices, and leads to slightly different final evaluation scores. Here we show the average and standard deviation of 5 runs with different random seeds to show the effect of random grouping. As we can see from Table 5, the selection of choices only lead to small changes in overall performance with standard deviation less than 1 for all runs.

D.2 Zero-shot Results for Flan-T5_{XL}

Since Flan-T5 is tuned with multi-choice format, it can also do element selection in zero-shot. However, as we can see from Table 6, while the model still gets some elements correct, it is much lower compared to the fine-tuned model, and 3-shot GPT 3.5/4. This is expected, since Flan-T5 is not tuned for HTML and coding related tasks.

Table 6: Zero-shot element selection results for Flan-T5_{XL} compared with the fine-tuned counterpart.

	Cross-Task	Cross-Website	Cross-Domain
Flan-T5 _{XL} Zero-Shot	10.8	7.8	11.7
Flan-T5 _{XL} Fine-Tuned	52.0	38.9	39.6

Table 7: Step Success Rate for all methods on the 50 tasks subsets we used to evaluate GPT-4. Numbers in parentheses are the results on the full test set (same as Table 2)

	Cross-Task	Cross-Website	Cross-Domain
Flan-T5 _B	43.3 (41.0)	25.3 (29.5)	28.1 (31.6)
Flan-T5 _L	48.1 (50.3)	30.8 (35.3)	27.6 (37.3)
Flan-T5 _{XL}	47.9 (52.0)	33.3 (38.9)	34.6 (39.6)
GPT-3.5	15.2 (17.4)	15.1 (16.2)	16.7 (18.6)
GPT-4	36.2	30.1	26.4

D.3 Results on the 50 task subsets

Due to budget constraint, we only run GPT-4 on 50 tasks for each setting. Here we show the step success rate results for other methods on the same 50 examples that GPT-4 is tested on. As we can see from Table 7, the results on the 50 tasks subsets are consistent with the results on the respective full test set, and the relative performance across methods and splits remains the same.

Table 8: Prompt for action prediction in MINDACT with GPT models.
Only part of the HTML snippet is shown here to save space.

Role	Content
system	You are a helpful assistant that is great at website design, navigation, and executing tasks for the user
user	<pre>``` <html> <div> <div> <a tock home page /> <button id=0 book a reservation. toggle open> Book a reservation </button> <button book a reservation. toggle open> </button> </div> <div> <select id=1 type> <option reservations true> Dine in </option> ... </html> ``` Based on the HTML webpage above, try to complete the following task: Task: Check for pickup restaurant available in Boston, NY on March 18, 5pm with just one guest Previous actions: None What should be the next action? Please select from the following choices (If the correct action is not in the page above, please select A. 'None of the above'):</pre> <ul style="list-style-type: none"> A. None of the above B. <button id=0 book a reservation. toggle open> Book a C. <select id=1 type> <option reservations true> Dine in </option> <option D. <div id=2> <p> Celebrating and supporting leading women shaking up
assistant	<p>Answer: C. Action: SELECT Value: Pickup</p> <pre>``` Answer: C. Action: SELECT Value: Pickup</pre>
user	<pre>``` <html> <div> <main main> <section tabpanel> <div> <ul tablist> <li tab heading level 3 search and> <li id=0 tab heading level 3 search and> Hotel <li tab heading level 3 search and> <li tab heading level 3 search and> <div tabpanel> <div id=1> <div> Dates* <button button clear dates /> </div> ... </html> ``` Based on the HTML webpage above, try to complete the following task: Task: Compare the fare types to book a 1-adult ticket from Springfield, IL to Austin, TX for April 29th 2023 Previous actions: [combobox] Enter your departing city, airport name, or airpor... -> TYPE: SPRINGFIELD [button] Springfield, IL, US (SPI) -> CLICK [combobox] Enter your destination city, airport name, or airp... -> TYPE: AUSTIN [button] Austin, TX, US (AUS) -> CLICK What should be the next action? Please select from the following choices (If the correct action is not in the page above, please select A. 'None of the above'):</pre> <ul style="list-style-type: none"> A. None of the above B. <li id=0 tab heading level 3 search and> Hotel C. <div id=1> <div> Dates* <button button clear dates D. <ul id=2> <a mobile tools> <a open united's tiktok
assistant	Answer: A.

Continued on next page

Table 8 – continued from previous page

Role	Content
user	<pre>``` <html> <div> <nav main menu> <div button> Car Sales </div> <div id=0> <div> <div> <div> Buy A Car </div> <div> Plan Your Purchase </div> </div> <div> <h4> Its Tax Refund Time. Treat Yourself to an Upgrade. </h4> <p> With a variety of options, invest your refund in what you really want - a quality, used vehicle from Enterprise. </p> ... </html></pre> ```
	<p>Based on the HTML webpage above, try to complete the following task: Task: Find a mini van at Brooklyn City from April 5th to April 8th for a 22 year old renter.</p> <p>Previous actions:</p> <p>[searchbox] Pick-up & Return Location (ZIP, City or Airport) (... -> TYPE: Brooklyn [option] Brooklyn, NY, US Select -> CLICK</p> <p>What should be the next action? Please select from the following choices (If the correct action is not in the page above, please select A. 'None of the above'):</p> <ul style="list-style-type: none">A. None of the aboveB. <div id=0> <div> <div> <div> Buy A Car </div> <div>C. <div id=1> Enterprise Fleet Management </div>D. <button id=2 selected pick-up date 03/19/2023> 19 <p>assistant Answer: D. Action: CLICK</p>