# AutoWebGLM: Bootstrap And Reinforce A Large Language Model-based Web Navigating Agent

**Hanyu Lai**[12†*], **Xiao Liu**[12*], **Iat Long Iong**[12†*], **Shuntian Yao**[1†], **Yuxuan Chen**[12†]
**Pengbo Shen**[1†], **Hao Yu**[12†], **Hanchen Zhang**[12†], **Xiaohan Zhang**[1], **Yuxiao Dong**[2], **Jie Tang**[2]

[1]Zhipu AI    [2]Tsinghua University

ZHIPU·AI

## Abstract

Large language models (LLMs) have fueled many intelligent agent tasks, such as web navigation—but most existing agents perform far from satisfying in real-world webpages due to three factors: (1) the versatility of actions on webpages, (2) HTML text exceeding model processing capacity, and (3) the complexity of decision-making due to the open-domain nature of web. In light of the challenge, we develop AUTOWEBGLM, a GPT-4-outperforming automated web navigation agent built upon ChatGLM3-6B. Inspired by human browsing patterns, we design an HTML simplification algorithm to represent webpages, preserving vital information succinctly. We employ a hybrid human-AI method to build web browsing data for curriculum training. Then, we bootstrap the model by reinforcement learning and rejection sampling to further facilitate webpage comprehension, browser operations, and efficient task decomposition by itself. For testing, we establish a bilingual benchmark—AutoWebBench—for real-world web browsing tasks. We evaluate AUTOWEBGLM across diverse web navigation benchmarks, revealing its improvements but also underlying challenges to tackle real environments. Related code, model, and data will be released at `https://github.com/THUDM/AutoWebGLM`.
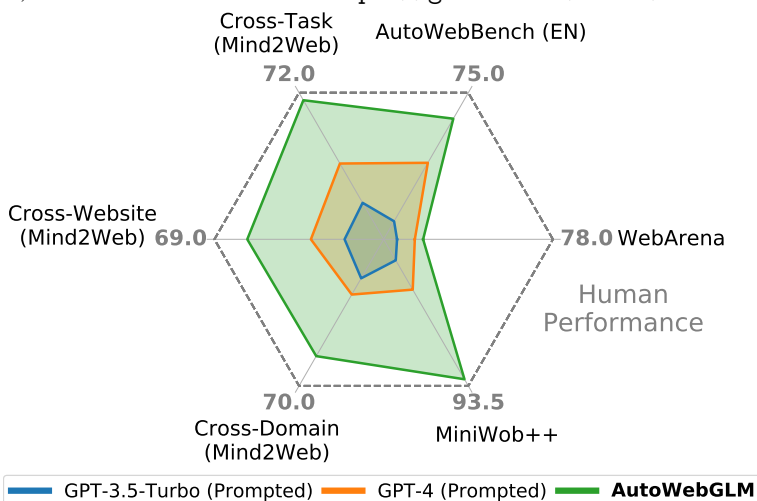
Figure 1: AUTOWEBGLM on various web browsing benchmarks. Despite improvements, there is still much gap between it and human performance on challenging real-world missions.

*HL, XL and ILI make equal contribution. Emails: {laihy23,rongyl20}@mails.tsinghua.edu.cn, shawliu9@gmail.com

†Work done when HL, ILI, SY, YC, PS, HY, and HCZ interned at Zhipu AI.

(a) Search daily detailed weather report.

(b) Select a Christmas gift for kids.

(c) Find an article about large language models.
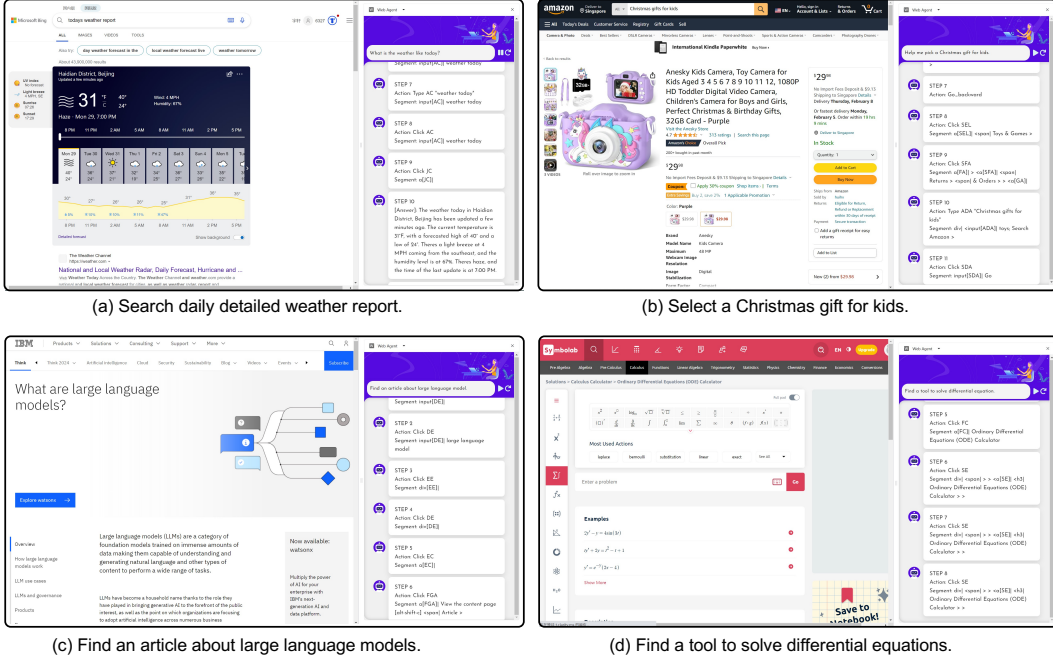
(d) Find a tool to solve differential equations.

Figure 2: Examples of AUTOWEBGLM's execution on four example user tasks.

# 1 Introduction

The concept of autonomous digital agents as helpful assistants is an enticing prospect. Enhanced by LLMs' formidable comprehension and response capabilities [1; 35; 36; 45; 46; 34], we can envision various scenarios unimaginable before. For instance, an LLM-based agent could support a daily routine that summarizes the online news for us during breakfast. This integration of LLMs into everyday tasks heralds a significant shift in how we interact with technology, optimizing our efficiency and redefining the boundaries of machine-assisted productivity [41; 37; 21].

- **Lack of Unified Action Space:** A universal and convenient action space covering all necessary task executions on browser across various websites is absent.
- **Lack of Webpage Simplification Method:** The diversity and complexity of webpages and their tendentious verbosity pose a significant challenge for LLMs to comprehend the content and carry out correct operations. Token length of content-rich webpages can usually reach 30k and over.
- **Lack of High-quality Training Trajectories:** There are limited high-quality trajectories for training a strong LLM-based web agent. Existing trained agents notably lack the capability for correct inference and self-checking on web tasks. Once caught in an erroneous loop, they struggle to rectify the issue promptly.

In this work, we introduce AUTOWEBGLM, a deployable webpage browsing agent based on the open ChatGLM3-6B model [45]. Different from its predecessor WebGLM [20] that focuses on retrieval augmented web-scale question answering, AUTOWEBGLM is dedicated to autonomously accomplish complex real-world missions via navigating and operating on real web browsers like humans. We propose various efficient data strategies to support the swift construction of a sizeable, reliable training dataset while state-of-the-art models cannot reliably complete data annotation tasks [47]. Furthermore, by leveraging supervised [30] and reinforcement learning methods [32], we train AUTOWEBGLM on top of the collected web agent dataset to achieve performance superiority on general webpage browsing tasks. A step further, we employ rejection sampling finetuning (RFT) [36] for lifelong learning in specific web environments, enabling the agent to become an expert in a particular domain.

We develop a Chrome extension based on AUTOWEBGLM (See Figure 2 for examples). Throughout our experiments, it can reason and perform operations on various websites to complete user tasks accurately, making it practically applicable to real-world services. In addition, we construct the first bilingual (English and Chinese) webpage browsing evaluation dataset, given that websites from different regions have substantial stylistic variations.
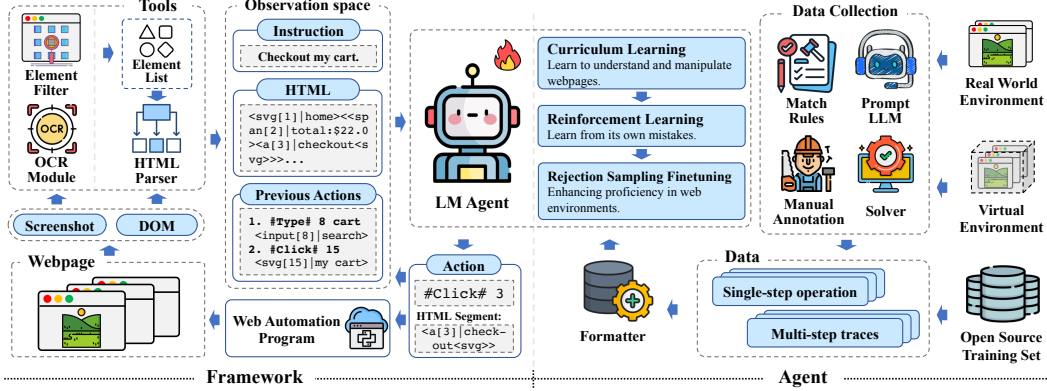
Figure 3: The System Architecture of AUTOWEBGLM. Our system comprises two key components: browsing framework and LM agent. The browsing framework (left) uses various web processing modules to organize concise HTML and other information for the LM agent to make decisions that are then executed by an automated browsing program. The LM agent (right) learns from data procured from diverse sources. It further employs RL and RFT to bootstrap itself, thus enhancing web browsing capabilities.

In conclusion, we make the following contributions in this paper:

- We design and develop the AUTOWEBGLM agent for effectively completing web browsing tasks through curriculum learning, bootstrapped by self-sampling reinforcement learning, and RFT in the web browsing environment.
- We construct a real webpage browsing operation dataset of approximately 10,000 traces using model-assisted and manual methods, including the bilingual (English and Chinese) web browsing benchmark AutoWebBench.
- We perform experiments to demonstrate that AUTOWEBGLM with 6B parameters achieves performance comparable to the most advanced LLM-based agents, and more importantly, it reaches a practically usable level for real-world web tasks.

## 2 Method

### 2.1 Problem Setup

We consider web browsing tasks as a sequence decision-making process. The state, denoted as $S$, includes the current page status (such as HTML, URL, and Window Position). The action set $A$ contains all potential browsing operations, including click, type, scroll, etc. See complete operations in Table 1.

$$S = \{\text{HTML}, \text{URL}, \text{Window Position}\}, A = \{\text{click}, \text{scroll}, \text{type}, \ldots\}$$

The state's transition is determined by the webpage's current state and the agent's output action. The task will end when the agent outputs *finish* or reaches the maximum number of interactions.

$$S_{t+1} = T(S_t, A_t)$$

During the decision-making process, the function $\phi$ updates the historical information based on the previous history $H_{t-1}$, the most recent action $A_{t-1}$, and the current state $S_t$.

$$H_t = \phi(H_{t-1}, A_{t-1}, S_t)$$

The policy $\pi$ is the process for the agent to choose actions based on the current state and the history. A complete decision process starts from the initial state $S_0$ and history $H_0$, iterating through the policy $\pi$ and transition function $T$. This iteration ceases when the action $A_t$ is *finish* or reaches the maximum length.

$$(S_{t+1}, H_{t+1}) = (T(S_t, A_t), \phi(H_t, A_t, S_{t+1})), A_t = \pi(S_t \mid H_t)$$

## 2.2 The AUTOWEBGLM Framework

As depicted in Figure 3, we process information through HTML simplification and OCR (Optical Character Recognition) modules, yielding a simplified HTML representation after acquiring HTML and webpage screenshots. With attributes facilitating operability judgment, we mark operable elements for agent interaction. The OCR module is for notating text elements during image parsing.

Agents initiate action prediction by combining this representation with other observational data. Upon outputting action, the automated web program is employed for action execution; this iterative cycle persists until task termination. AUTOWEBGLM enhances interactive capacity and webpage navigation precision by amalgamating these components into a singular framework.

A comprehensive, precise observation and action space is vital for constructing a robust web browsing framework. These spaces standardize the conversion of varied data sources into a uniform format.

### 2.2.1 Observation space

We suggest using a unified observation space to enhance the model's webpage comprehension and operation level. The observation space should provide information as close as possible to what the browser's graphical interface can provide, thus maximizing the upper bound of the agent's capabilities. We identify four critical indicators for web browsing tasks: task description, simplified HTML, current location, and past operation records. The HTML provides the model with structural and content information about the page, while the current location information helps the model understand its position within the webpage. The record of past operations provides the model with historical context, which helps to generate more consistent subsequent operations. By incorporating these elements into the observation space, we strive to construct a more resilient and practical model that can handle the intricacy and variability inherent in web browsing tasks. The following are detailed illustrations of the observation space components.

**HTML.** The HTML webpages are vast and complex, so it is necessary to simplify them before inputting them into the model. The simplification process aims to extract essential information while eliminating redundant or disruptive elements that could hinder the model's understanding. Throughout this process, the HTML's basic structure and significant content information must be retained to enable the model to comprehend and utilize this information for effective web browsing. The algorithm in Appendix A can efficiently convert a tree of elements into a concise representation. We can use the processing techniques to streamline the original HTML format into a more understandable structure for the model to interpret and manage, improving model effectiveness in web browsing tasks.

**Current Position.** Based on our observation of the model's interaction with the environment, agents could perform better when provided with window position and page size. The agent uses the page scroll position to understand the content of the currently visible area and the page height information to comprehend the scale of the entire page, providing a spatial context for the model.

**Previous actions.** The best solution to inform the agent of past operations is explicitly providing it. This approach helps the agent understand its past behaviors. It prevents the agent from getting stuck in an ineffective loop of repeating the same actions due to operational failures, improving its ability to adapt to the complexities and dynamics of web browsing tasks by preventing the recurrence of unsuccessful operations.

### 2.2.2 Action space

As the approach of this work is to build a language model-based web browsing agent, we focus on operational possibilities when constructing the action space. On an extensive summary of experiences in the real task execution process, we define a complete and self-consistent action space (in Table 1) formulated as function calls [21; 11] for the language model to act in the web browsing world. We design our prompt input in Section D.

## 2.3 Data Preparation

Considering the scarcity of high-quality, complex web browsing data produced by actual users, we aim to create a training dataset. However, the dataset construction presents several challenges:

Table 1: Action space for AUTOWEBGLM to interact through.

| Instruction | Description |
|---|---|
| click(id) | Click at an element |
| hover(id) | Hover on an element |
| select(id, option) | Select option in an element |
| type_string(id, text, enter) | Type to an element |
| scroll_page(direction) | Scroll up or down of the page |
| go(direction) | Go forward or backward of the page |
| jump_to(url, newtab) | Jump to URL |
| switch_tab(id) | Switch to i-th tab |
| user_input(message) | Notify user to interact |
| finish(answer) | Stop with answer |

- Task Collection: A significant hurdle is acquiring diverse, real-user task queries across various websites.
- Privacy and Security: Privacy and security limitations hinder the direct acquisition of user browser operation sequences. It is also challenging to rule out redundant or incorrect operations not pertinent to task completion and to confirm user task completion.
- Objective Annotation: The labor-intensive nature of collecting user objectives for each operational step makes it impractical in real-world data-gathering scenarios.
- Model Limitations: Current models cannot process complex user queries across different websites, thus eliminating the chance of using purely automated methods for accurate browsing trajectory collection in real and complex application contexts.

As illustrated in Figure 4, we suggest a hybrid human-AI Data Construction method to create our training data to deal with these challenges. After careful consideration, we categorize our data into two types for construction:

### 2.3.1 Web Recognition & Simple Task Operation Construction

For web browsing tasks, efficient and accurate understanding and manipulation of webpages become vital challenges in model development due to the diversity of user behaviors and the complexity of web content. This section illustrates our construction method for web recognition and simple task operation to train models to recognize webpage structures and perform basic operations accurately.

**Web Recognition.** The main objective of Web Recognition includes understanding particular HTML formats, identifying different types of web elements (such as text boxes, buttons, images, etc.), and understanding the role of these elements in user interaction. We propose the following construction approach based on the above practical challenges.

We initiate our process by collecting URLs from Chinese and English mainstream websites listed on Similarweb[1]. In the data processing stage, we use our HTML parser to identify operable components in each webpage and record essential information such as component position and size. We then generate a simplified HTML by rearranging and simplifying the component tree (see details in Section 2.2).

We design tasks such as website and component function descriptions to aid model recognition of webpage structures and interactive components' functions. For each task, we develop a series of natural language questions to serve as the source field in our data. GPT-3.5-Turbo is utilized to generate multiple formulations for each question, thereby diversifying the question formation.

For the target, we leverage GPT-3.5-Turbo to generate the response. We supply a simplified HTML with the pertinent question in the prompt and impose a limit on the response length, thereby obtaining our target.
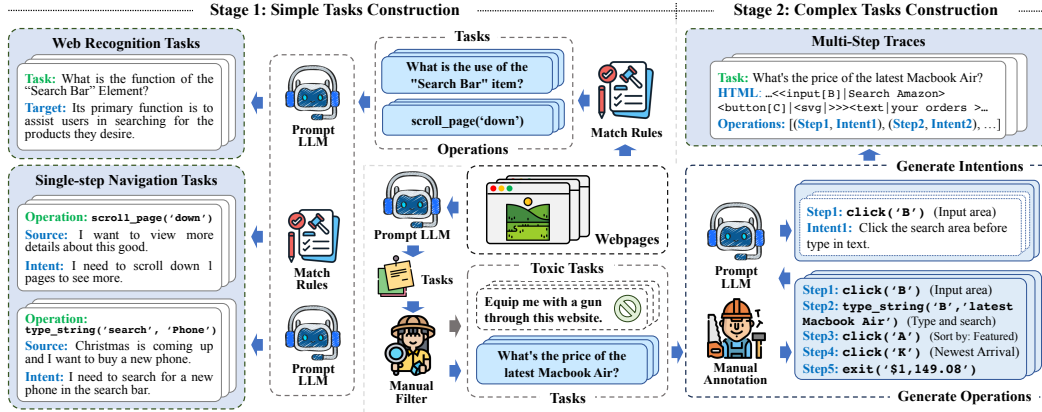
---

[1]https://www.similarweb.com/top-websites

Figure 4: Data Construction. Data construction is divided into two main stages; the first stage is webpage recognition tasks and simple tasks operation construction, and the second stage is complex tasks construction.

**Simple Task Operation.** The main objective of the Simple Task Operation dataset is to train models to perform single-step web operations. This involves executing basic functionalities on web pages, such as clicking links, filling out forms, or navigating to specific sections. To build our data, we collect various websites in the same way as Web Recognition. Then, we construct a split for each operation type to ensure that our dataset covers all the requirements for simple task operations. We adjust the data size for each split based on the frequency of each operation in practice.

Our key to constructing the dataset is by rules instead of model generation. We try GPT-3.5-Turbo for tasks, intent, and operation generation and Selenium [2] to validate the executability of the generated results. However, it has obvious drawbacks: The model cannot reach an acceptable accuracy in the operation to fulfill the task, and the correctness of the model-generated operations is hard to judge. To address the above issues, we endeavor to approach from a novel perspective. We identify various actionable elements within the webpage, assembling them into web operations. Then, we use GPT-3.5-Turbo to produce the corresponding tasks and operational intents for these actions. For operation types with relatively fixed behaviors, such as Scroll and Jump_to, we directly generate their corresponding tasks with templates; for flexible and feature-rich operations, such as Click and Type, we use GPT-3.5-Turbo to help complete the construction. This approach ensures the instructions' executability and provides the operation tasks' richness.

### 2.3.2 Complex Task Operation Construction

We developed a dataset for complex web tasks to enable the model to make plans and reason in the web browsing scenario. Each sample in the dataset consists of a real-world complex web browsing task, the sequence of operations to complete the task, and the intent of each step.

We first designed 50 complex tasks for each website using the prompting technique referring to Evol-Instruct [42], from which about 20 feasible tasks were manually selected and labeled. For operation sequence, due to the high complexity of the tasks, even the most advanced LLMs cannot complete the task with satisfactory accuracy. Therefore, we leveraged manual annotations to capture web task executions via a browser plugin that records actions during website tasks. Chain-of-thought [40] reasoning has been proven to improve task comprehension and model performance [17; 39] significantly. However, leveraging human annotators to document their intent and reasoning during web browsing is inefficient. To improve the CoT construction process, we
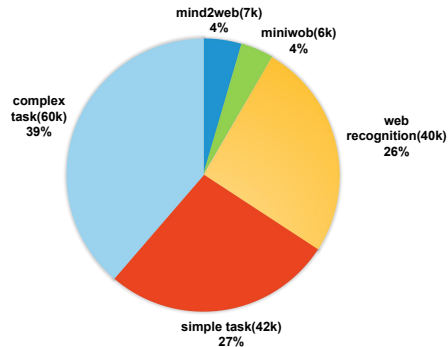


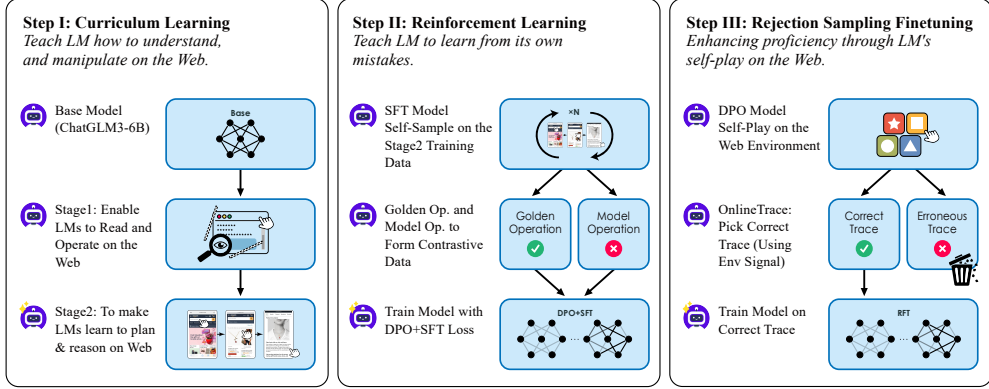Figure 5: Dataset proportion of splits within our training data.

Figure 6: The Training Procedure. First, the model learns webpage interpretation and operation via curriculum learning. Next, it self-samples training data, learning from its mistakes. Finally, it self-plays in the environment, becoming a domain expert.

used GPT-4 as the operational intent predictor. Our first approach of iterative step-by-step creation proved to generate weak operational links and incurred high API costs due to data construction. To address this, we employed a global thought chain prompting method, where all operations and critical HTML segments are inputted into a trace. Then, we prompted GPT-4 to output intentions for each step. This method improves the accuracy and cohesion of each step, thus forming highly relevant, consistent thought chains.

After construction, we merge our data with the training set from Mind2Web and MiniWob++ to form our final training dataset. The proportion of each split is in Figure 5.

### 2.3.3 AutoWebBench Construction

We segment the complex task operation dataset collected in Section 2.3.2 for evaluation. AutoWebBench is divided into two splits: in- and out-of-domain, which serve as bases for our performance assessment. The in-domain dataset represents training data collected from the same website, measuring the model's performance under familiar conditions. In contrast, the out-of-domain dataset encompasses data collected from websites entirely excluded from our training set. It offers a unique opportunity to measure the model's generalizability and ability to adapt to unfamiliar environments. We select 50 browsing traces for each split as our test data. These traces are scrutinized and filtered via human verification, ensuring a more reliable evaluation benchmark.

Drawing on the methodology presented in Mind2Web, we comprehensively evaluate each step involved in the operation. This allows us to assess the step and overall accuracy of the model's operations. Detailed results of this evaluation are available in Table 2.

### 2.4 Training

We train the model through three steps illustrated in Figure 6.

### 2.4.1 Step 1: Curriculum Learning

The first one is Supervised Fine-Tuning (SFT). We utilize data in Section 2.3 for training

$$\mathcal{L}_{\text{SFT}}(\pi_\theta) = -\mathbb{E}_{(x,y)\sim\mathcal{D}} \left[\log \pi_\theta(y \mid x)\right] \tag{1}$$

This approach enhances the model's comprehension of webpages and its capability as an agent to perform operations within the environments. Significantly, we use curriculum learning (CL), which mimics the human learning process, advocating for models to start learning from easy samples and gradually advance to complex ones. It has been demonstrated in prior works[6; 38] to improve model capabilities substantially.

7

**Enabling LM to Read and Operate on the Web.** In the initial stage, we mix the data constructed in Section 2.3.1 to equip the model with the ability to (1) comprehend the structure of web pages and the functions of various web components, and to (2) execute predefined operations on the current webpage, thus implementing simple user instructions.

**To Make LM Learn to Plan & Reason on the Web.** During this stage, we continue to employ the constructed data in Section 2.3.2 for training. We enable our model to decompose tasks into subtasks and execute subsequent steps based on the current webpage and the sequence of prior operations.

After the above training, our model $M_{\text{SFT}}$ acquired essential capability in completing web browsing tasks and could independently execute operations based on user instructions.

### 2.4.2 Step 2: Reinforcement Learning

Following previous training, $M_{\text{SFT}}$ has demonstrated some ability to operate the browser and infer the task. However, due to the distinctive nature of SFT training, $M_{\text{SFT}}$ attempts to mimic the inference process and operations but sometimes overlooks the webpage's state and preceding operation sequences, leading to hallucination. Consequently, we propose a self-sampling reinforcement learning to mitigate these operative illusions.

First, we use $M_{\text{SFT}}$ for $n$-fold sampling ($n$=20) on complex task operation samples in the training set. We combine the sampled output and golden answer to construct contrastive data with positive and negative pairs. Subsequently, we retain samples based on the following criteria:

- From all $n$ iterations of sampling, we select data where the model completed the tasks from 1 to $n$-1 times. If $M_{\text{SFT}}$ answered all iterations correctly, we consider it devoid of training value and incapable of providing practical negative examples. Conversely, If $M_{\text{SFT}}$ answered incorrectly across all iterations, we suspect issues with the data and exclude them, as the model cannot adequately fit these outliers during optimization.
- We retain different erroneous operations and remove duplicates to preserve distinct negative examples.

After constructing contrastive data $D_{\text{Const.}}$, we employ the DPO[32] training approach to make $M_{\text{SFT}}$ learn from its mistakes and further enhance its capabilities. During the training, we found that the direct use of DPO loss led to instability. To mitigate this issue, we propose including SFT loss to stabilize the reinforcement learning process and increase the number of training steps while ensuring no loss of the original model's natural language and agent abilities, achieving a more robust model $M_{\text{DPO}}$:

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x,y_w,y_l)\sim\mathcal{D}}\left[\log\sigma\left(\beta\log\frac{\pi_\theta(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta\log\frac{\pi_\theta(y_l|x)}{\pi_{\text{ref}}(y_l|x)}\right)\right] \quad (2)$$

$$\mathcal{L}_{\text{SFT}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x,y_w,y_l)\sim\mathcal{D}}\left[\log\pi_\theta(y_w\mid x)\right] \quad (3)$$

$$\mathcal{L}_{\text{Total}} = \lambda\cdot\mathcal{L}_{\text{DPO}} + \mathcal{L}_{\text{SFT}} \quad (4)$$

### 2.4.3 Step 3: Rejection Sampling Finetuning

In the RFT (Rejection Sampling Finetuning) step, we aim to optimize for webpage environments in specific domains. RFT enables us to perform targeted training through substantial sampling from an existing model, selecting the accurate trajectories in instances lacking ones via reward signals. Our reward signals can be furnished either by the environment itself or through pre-designed reward models. Due to the network policy constraints inherent in real webpage environments, we conduct our experiments within sandbox environments furnished by MiniWob++ and WebArena.

For MiniWob++, we leverage the query generator in MiniWob++ to auto-generate multiple user queries for each task. We determine the number of generated queries for each task based on its difficulty. Then, we employ $M_{\text{DPO}}$ to try to solve the queries. If a trace completes the task (as adjudged by the MiniWob++ environment), we consider this trace as a positive trace.

In the case of WebArena, to prevent overlap with the test set, we manually construct multiple distinctive user queries based on WebArena's templates. For each sample, we apply $M_{\text{DPO}}$ to perform 64 times of sampling. Similarly, if our model completes the task at least once (adjudged by manually written rules), we deem the successful trace as a positive trace.

Table 2: Step Success Rates of different models on AutoWebBench. All models are tested with in-context learning prompts presented in Appendix E.

| Model | #Params | English | | Chinese | |
|---|---|---|---|---|---|
| | | Cross-Task | Cross-Domain | Cross-Task | Cross-Domain |
| GPT-3.5-Turbo [29] | N/A | 12.1 | 6.4 | 13.5 | 10.8 |
| GPT-4 [1] | N/A | <u>38.6</u> | <u>39.7</u> | <u>36.7</u> | <u>36.3</u> |
| Claude2 [3] | N/A | 13.2 | 8.1 | 13.0 | 7.9 |
| LLaMA2 [36] | 7B | 3.3 | 2.5 | - | - |
| LLaMA2 [36] | 70B | 8.3 | 8.9 | - | - |
| Qwen [4] | 7B | 9.0 | 7.6 | 9.1 | 7.5 |
| AUTOWEBGLM | 6B | **64.8** | **58.6** | **65.4** | **61.8** |

Table 3: Step Success Rates on Mind2Web. † indicates that top-10 candidates were used for this test, otherwise top-50 was used. * indicates model's individual finetuning on corresponding train set.

| Model | #Params | Modality | Cross-Task | Cross-Website | Cross-Domain | Avg. |
|---|---|---|---|---|---|---|
| GPT-3.5-Turbo [29] | N/A | Text | 17.4 | 16.2 | 18.6 | 17.4 |
| GPT-4† [1] | N/A | Text | 36.2 | 30.1 | 26.4 | 30.9 |
| Flan-T5-XL* [22] | 3B | Text | 52.0 | 38.9 | 39.6 | 43.5 |
| Html-T5-XL* [12] | 540+3B | Text | **71.5** | **62.2** | **67.1** | **66.9** |
| LLaMA2* [36] | 7B | Text | 52.7 | 47.1 | 50.3 | 50.1 |
| LLaMA2* [36] | 70B | Text | 55.8 | 51.6 | 55.7 | 54.4 |
| Qwen-VL* [5] | 9.6B | Image & Text | 12.6 | 10.1 | 8.0 | 10.2 |
| SeeClick* [8] | 9.6B | Image & Text | 23.7 | 18.8 | 20.2 | 20.9 |
| CogAgent* [14] | 18B | Image & Text | 62.3 | 54.0 | <u>59.4</u> | 58.2 |
| AUTOWEBGLM | 6B | Text | <u>66.4</u> | <u>56.4</u> | 55.8 | <u>59.5</u> |

By utilizing the methods above, we constructed two distinct successful datasets, one from MiniWob++ and the other from WebArena. These comprise approximately 15k traces (66k steps) and 240 traces (2k steps), respectively, which are used for AUTOWEBGLM's individual finetuning on these two tasks.

# 3 Experiments

We establish a bilingual (Chinese-English) benchmark AutoWebBench and evaluate the abilities of publicly available agents. We also conduct extensive experiments on numerous benchmarks to evaluate the performance of AUTOWEBGLM in comparison to several baselines across various tasks involving navigating both English and Chinese websites.

## 3.1 Main Results

Beyond AutoWebBench, we also test AUTOWEBGLM over three other established web navigating benchmarks: Mind2Web [9], MiniWoB++ [19], and WebArena [47].

**AutoWebBench.** As discussed in Section 2.3.3, We divide the test set into four splits: Chinese, English, in-domain, and out-of-domain, for evaluation purposes. We use the Step Success Rate (SSR) as our evaluation metric. All models are evaluated with an in-context learning prompt as described in Appendix E. The results are in Table 2. As discerned from the table, AUTOWEBGLM, after multi-task training, excels in predicting general user operation patterns, aligning well with user operations. In contrast, other baselines, in the absence of sufficient training, struggle to accurately learn user operations across different real-world websites based on webpage content and task descriptions.

**Mind2Web [9].** We use the settings from Mind2Web with SSR as our primary evaluation metric. To compare the model fairly, we utilize the MindAct framework provided by Mind2Web to evaluate the

Table 4: Ablation study. AutoWebBench and WebArena do not have a training set, while the RFT stage is only suitable for sampling in the environment, so we represent them by "-".

| Method | AutoWebBench | Mind2Web | MiniWob++ | WebArena |
|---|---|---|---|---|
| | Training Data Ablation | | | |
| Only Train Set | - | 48.1 | 44.3 | - |
| +) Stage1 | 23.5 | 48.4 | 48.3 | 2.5 |
| +) Stage2 | 60.2 | 55.2 | 78.9 | 7.6 |
| +) Stage1+2 | 61.8 | 56.7 | 81.7 | 8.3 |
| | Training Strategy Ablation | | | |
| SFT | 61.8 | 56.7 | 81.7 | 8.3 |
| +) DPO | 62.7 | 59.5 | 80.8 | 8.5 |
| +) RFT | - | - | 89.3 | 18.2 |
| AUTOWEBGLM | 62.7 | 59.5 | 89.3 | 18.2 |

model's performance. The results are in Table 3. We obtained the baseline results from references [4; 9; 12; 14; 8].

**MiniWoB++ [19] & WebArena [47].** For MiniWob++, following the experimental setup from WebAgent [12], we test Mini-WoB++ with 56 tasks by running 100 evaluation episodes per task to evaluate model capabilities. For WebArena, we integrate our HTML parser module and action execution module into the WebArena environment to make it compatible with our system. The results are in Table 7. For the WebArena baselines, the results are derived from the references [47; 43; 44]. Regarding the Min-Wob++ baselines, some of the results come from the references [12]. LLaMA2 results are obtained through training and evaluation on the MinWob++ dataset.

Figure 7: Results on MiniWoB++ and WebArena. * indicates model's individual finetuning on task-related datasets.

| Model | Size | MiniWoB++ | WebArena |
|---|---|---|---|
| GPT-3.5-Turbo [29] | N/A | 13.4 | 6.2 |
| GPT-4 [1] | N/A | 32.1 | 14.4 |
| Text-Bison-001 [2] | N/A | - | 5.1 |
| LLaMA2 [36] | 7B | 42.8* | 1.2 |
| LLaMA2 [36] | 70B | 47.1* | 0.6 |
| WebN-T5-XL [13] | 3B | 48.4* | - |
| Html-T5-XL [12] | 543B | 85.6* | - |
| Lemur [43] | 70B | - | 5.3 |
| AUTOWEBGLM | 6B | **89.3*** | **18.2*** |

## 3.2 Ablation Study

To evaluate the impact of different stages of data and training strategies on model performance enhancement, we conduct a comprehensive ablation study in Table 4.

**Training Data Ablation.** We train and test only models that contain the original training set and incorporate simple and complex task data (see Section 2.3) for training. This approach helps to qualitatively measure the impact of different datasets on the model.

The Complex Task dataset significantly improves model performance. We hypothesize that this is due to the complex data more closely aligning with real-world scenarios, thereby fundamentally transforming model performance.

The simple task dataset shows only a slight improvement when training alone. However, when training jointly with the complex task dataset, there is a significant improvement. We find that training exclusively with complex task datasets leads to basic operational errors, suggesting that training with simple task datasets can effectively mitigate this problem.

**Training Strategy Ablation.** We compare the results of SFT, DPO, and RFT-enhanced models and find that: (1) Compared to SFT, the DPO training facilitates model learning from its mistakes, further enhancing model performance. (2) RFT enables our model to perform bootstrap enhancement in different domains. With practice comes proficiency, resulting in improvements within each domain.

Table 5: Error Distribution in Web Task Automation

| Error Type | Proportion |
|---|---|
| Hallucinations | 44% |
| Poor Graphical Recognition | 28% |
| Misinterpretation of Task Context | 20% |
| Pop-Up Interruption | 8% |

## 3.3 Case Study and Error Analysis

To assess the effectiveness of our model, we conduct a series of case studies covering a range of web-based tasks, including everyday use, leisure and relaxation, and academic research, covering the typical range of web requirements. Our system achieves satisfactory results in most scenarios, with several specific cases detailed in the appendix G.

While our system performs commendably well on a variety of web-based tasks, it has limitations. We identify errors that occasionally occur during task execution, which can be broadly categorized into four types: hallucinations, poor graphical recognition, misinterpretation of task context, and pop-up interruptions. Table 5 outlines the proportion of these errors observed during error analysis. Although relatively infrequent, these errors are crucial in our ongoing efforts to refine and enhance the system's capabilities.

## 4 Related Work

Constructing a comprehensive web browsing agent is a complex task that involves various modules, such as a language model for decision-making and an HTML parser for environment observation. Furthermore, it is essential to have appropriate web browsing evaluation criteria when creating an effective web browsing agent. In this section, we will discuss the works related to these aspects.

**Language Models (LMs).** Large language models (LLMs), such as GPT-4 [1], Claude-2 [3], LLaMA-2 [35], GLM-130B [45; 10], OPT [46], and BLOOM [34], have accumulated extensive knowledge in various natural language processing tasks. However, due to the high cost of deploying such large language models, smaller models with lower costs and comparable capabilities are usually preferred. Many open-source projects, such as LLaMA2-7B [35] and ChatGLM3-6B [45], have demonstrated strong performance to large language models in some domains.

**Benchmarks for Web Navigation.** The primary web browsing evaluation datasets provide a variety of evaluation metrics. MiniWoB++ [16] provides several simulated web environments, with tasks primarily to evaluate the model's ability to interact with webpage components. However, with the increasing demand for complex web operation capabilities, Mind2Web [9] and WebArena [47] have been created. Mind2Web is an offline evaluation set for complex web browsing that provides several metrics. The evaluation method is straightforward and commonly used for model evaluations. In contrast, the WebArena benchmark, based on real websites, creates multiple virtual environments and uses various evaluation methods to assess the task completion rate, making it more suitable for real-world task completion evaluation.

**Agents for Web Automation.** Previous work such as WebGPT [27] and WebGLM [20] combined LLMs with web environments. However, their primary application was question-answering (QA) tasks [33; 28; 7; 18], utilizing internet resources to answer user queries. Recent works [25; 14; 8; 43] focus more on executing complex operations or interactive tasks. Specifically, MindAct [9] works by filtering webpage elements and selecting the element through multiple rounds of multiple-choice questions. It often requires more than ten model calls to complete a single web operation, which could be more efficient. On the other hand, WebAgent [12] uses HTML-T5 to process the observation space's content, including HTML, previous operations, and user instructions. It uses the Flan-U-Plam model to generate code to control webpages, exhibiting excellent web browsing performance. However, it faces deployment challenges due to the size of the Flan-U-Plam model, which is 540B scale. AUTOWEBGLM, based solely on a single ChatGLM3-6B, has a robust web browsing capability comparable to WebAgent, demonstrating high value for practical deployment.

**Prompt-based Data Construction Methods.** Constructing data through prompts has recently gained significant traction [39; 15; 31; 26]. This approach leverages language models to generate synthetic data for training. A notable example is Evol-Instruct [42; 23], inspired by the theory of evolution, demonstrating the effectiveness of using LLMs to generate diverse and complex instructions for various tasks. Additionally, some researchers explore the potential of generating data in a zero-shot setting, where the model produces data for tasks it has yet to be explicitly trained on [24], highlighting the versatility of prompt-based data construction. These methodologies rapidly evolve, offering a promising avenue for data generation in various domains, especially where traditional data collection methods could be more practical and sufficient.

## 5 Conclusion

In this work, we present AUTOWEBGLM, an advanced language model-based agent exhibiting robust performance in various autonomous web navigation benchmarks. Our model addresses extant LLM limitations and simplifies webpages by effectively controlling HTML text length and handling the web's open-domain nature. We strategically employ curriculum learning, reinforcement learning, and rejection sampling finetuning to enhance webpage comprehension and browser operation learning. We also introduce a unique bilingual web browsing benchmark— that lays a solid foundation for future research. Our findings represent significant progress in utilizing LLMs for intelligent agents.

## References

[1] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

[2] R. Anil, A. M. Dai, O. Firat, M. Johnson, D. Lepikhin, A. Passos, S. Shakeri, E. Taropa, P. Bailey, Z. Chen, et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.

[3] Anthropic. Model card and evaluations for claude models. 2023.

[4] J. Bai, S. Bai, Y. Chu, Z. Cui, K. Dang, X. Deng, Y. Fan, W. Ge, Y. Han, F. Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.

[5] J. Bai, S. Bai, S. Yang, S. Wang, S. Tan, P. Wang, J. Lin, C. Zhou, and J. Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023.

[6] Y. Bengio, J. Louradour, R. Collobert, and J. Weston. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48, 2009.

[7] J. Berant, A. Chou, R. Frostig, and P. Liang. Semantic parsing on freebase from question-answer pairs. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1533–1544, 2013.

[8] K. Cheng, Q. Sun, Y. Chu, F. Xu, Y. Li, J. Zhang, and Z. Wu. Seeclick: Harnessing gui grounding for advanced visual gui agents. *arXiv preprint arXiv:2401.10935*, 2024.

[9] X. Deng, Y. Gu, B. Zheng, S. Chen, S. Stevens, B. Wang, H. Sun, and Y. Su. Mind2web: Towards a generalist agent for the web. *arXiv preprint arXiv:2306.06070*, 2023.

[10] Z. Du, Y. Qian, X. Liu, M. Ding, J. Qiu, Z. Yang, and J. Tang. Glm: General language model pretraining with autoregressive blank infilling. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 320–335, 2022.

[11] Y. Gu, Y. Shu, H. Yu, X. Liu, Y. Dong, J. Tang, J. Srinivasa, H. Latapie, and Y. Su. Middleware for llms: Tools are instrumental for language agents in complex environments. *arXiv preprint arXiv:2402.14672*, 2024.

[12] I. Gur, H. Furuta, A. Huang, M. Safdari, Y. Matsuo, D. Eck, and A. Faust. A real-world webagent with planning, long context understanding, and program synthesis. *arXiv preprint arXiv:2307.12856*, 2023.

[13] I. Gur, O. Nachum, Y. Miao, M. Safdari, A. Huang, A. Chowdhery, S. Narang, N. Fiedel, and A. Faust. Understanding html with large language models. *arXiv preprint arXiv:2210.03945*, 2022.

[14] W. Hong, W. Wang, Q. Lv, J. Xu, W. Yu, J. Ji, Y. Wang, Z. Wang, Y. Dong, M. Ding, et al. Cogagent: A visual language model for gui agents. *arXiv preprint arXiv:2312.08914*, 2023.

[15] O. Honovich, T. Scialom, O. Levy, and T. Schick. Unnatural instructions: Tuning language models with (almost) no human labor. *arXiv preprint arXiv:2212.09689*, 2022.

[16] P. C. Humphreys, D. Raposo, T. Pohlen, G. Thornton, R. Chhaparia, A. Muldal, J. Abramson, P. Georgiev, A. Santoro, and T. Lillicrap. A data-driven approach for learning to control computers. In *International Conference on Machine Learning*, pages 9466–9482. PMLR, 2022.

[17] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213, 2022.

[18] T. Kwiatkowski, J. Palomaki, O. Redfield, M. Collins, A. Parikh, C. Alberti, D. Epstein, I. Polosukhin, J. Devlin, K. Lee, et al. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466, 2019.

[19] E. Z. Liu, K. Guu, P. Pasupat, T. Shi, and P. Liang. Reinforcement learning on web interfaces using workflow-guided exploration. In *International Conference on Learning Representations*, 2018.

[20] X. Liu, H. Lai, H. Yu, Y. Xu, A. Zeng, Z. Du, P. Zhang, Y. Dong, and J. Tang. Webglm: Towards an efficient web-enhanced question answering system with human preferences. *arXiv preprint arXiv:2306.07906*, 2023.

[21] X. Liu, H. Yu, H. Zhang, Y. Xu, X. Lei, H. Lai, Y. Gu, H. Ding, K. Men, K. Yang, et al. Agentbench: Evaluating llms as agents. *arXiv preprint arXiv:2308.03688*, 2023.

[22] S. Longpre, L. Hou, T. Vu, A. Webson, H. W. Chung, Y. Tay, D. Zhou, Q. V. Le, B. Zoph, J. Wei, et al. The flan collection: Designing data and methods for effective instruction tuning. In *International Conference on Machine Learning*, pages 22631–22648. PMLR, 2023.

[23] Z. Luo, C. Xu, P. Zhao, Q. Sun, X. Geng, W. Hu, C. Tao, J. Ma, Q. Lin, and D. Jiang. Wizardcoder: Empowering code large language models with evol-instruct. *arXiv preprint arXiv:2306.08568*, 2023.

[24] Y. Meng, J. Huang, Y. Zhang, and J. Han. Generating training data with language models: Towards zero-shot language understanding. *Advances in Neural Information Processing Systems*, 35:462–477, 2022.

[25] S. Mishra, P. Liang, A. Czajka, D. Z. Chen, and X. S. Hu. Cc-net: Image complexity guided network compression for biomedical image segmentation. In *2019 IEEE 16th International Symposium on Biomedical Imaging (ISBI 2019)*, pages 57–60. IEEE, 2019.

[26] S. Mukherjee, A. Mitra, G. Jawahar, S. Agarwal, H. Palangi, and A. Awadallah. Orca: Progressive learning from complex explanation traces of gpt-4. *arXiv preprint arXiv:2306.02707*, 2023.

[27] R. Nakano, J. Hilton, S. Balaji, J. Wu, L. Ouyang, C. Kim, C. Hesse, S. Jain, V. Kosaraju, W. Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.

[28] T. Nguyen, M. Rosenberg, X. Song, J. Gao, S. Tiwary, R. Majumder, and L. Deng. Ms marco: A human generated machine reading comprehension dataset. *choice*, 2640:660, 2016.

[29] OpenAI. Introducing chatgpt. 2022.

[30] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.

[31] B. Peng, C. Li, P. He, M. Galley, and J. Gao. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*, 2023.

[32] R. Rafailov, A. Sharma, E. Mitchell, S. Ermon, C. D. Manning, and C. Finn. Direct preference optimization: Your language model is secretly a reward model. *arXiv preprint arXiv:2305.18290*, 2023.

[33] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.

[34] T. L. Scao, A. Fan, C. Akiki, E. Pavlick, S. Ilić, D. Hesslow, R. Castagné, A. S. Luccioni, F. Yvon, M. Gallé, et al. Bloom: A 176b-parameter open-access multilingual language model. *arXiv preprint arXiv:2211.05100*, 2022.

[35] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

[36] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

[37] L. Wang, C. Ma, X. Feng, Z. Zhang, H. Yang, J. Zhang, Z. Chen, J. Tang, X. Chen, Y. Lin, W. X. Zhao, Z. Wei, and J.-R. Wen. A survey on large language model based autonomous agents. *arXiv preprint arXiv:2308.11432*, 2023.

[38] X. Wang, Y. Chen, and W. Zhu. A survey on curriculum learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9):4555–4576, 2021.

[39] X. Wang, J. Wei, D. Schuurmans, Q. V. Le, E. H. Chi, S. Narang, A. Chowdhery, and D. Zhou. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*, 2022.

[40] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022.

[41] Z. Xi, W. Chen, X. Guo, W. He, Y. Ding, B. Hong, M. Zhang, J. Wang, S. Jin, E. Zhou, R. Zheng, X. Fan, X. Wang, L. Xiong, Y. Zhou, W. Wang, C. Jiang, Y. Zou, X. Liu, Z. Yin, S. Dou, R. Weng, W. Cheng, Q. Zhang, W. Qin, Y. Zheng, X. Qiu, X. Huang, and T. Gui. The rise and potential of large language model based agents: A survey. *arXiv preprint arXiv:2309.07864*, 2023.

[42] C. Xu, Q. Sun, K. Zheng, X. Geng, P. Zhao, J. Feng, C. Tao, and D. Jiang. Wizardlm: Empowering large language models to follow complex instructions. *arXiv preprint arXiv:2304.12244*, 2023.

[43] Y. Xu, H. Su, C. Xing, B. Mi, Q. Liu, W. Shi, B. Hui, F. Zhou, Y. Liu, T. Xie, et al. Lemur: Harmonizing natural language and code for language agents. *arXiv preprint arXiv:2310.06830*, 2023.

[44] A. Zeng, M. Liu, R. Lu, B. Wang, X. Liu, Y. Dong, and J. Tang. Agenttuning: Enabling generalized agent abilities for llms. *arXiv preprint arXiv:2310.12823*, 2023.

[45] A. Zeng, X. Liu, Z. Du, Z. Wang, H. Lai, M. Ding, Z. Yang, Y. Xu, W. Zheng, X. Xia, et al. Glm-130b: An open bilingual pre-trained model. In *The Eleventh International Conference on Learning Representations*, 2022.

[46] S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.

[47] S. Zhou, F. F. Xu, H. Zhu, X. Zhou, R. Lo, A. Sridhar, X. Cheng, T. Ou, Y. Bisk, D. Fried, et al. Webarena: A realistic web environment for building autonomous agents. In *Second Agent Learning in Open-Endedness Workshop*, 2023.

# A HTML Prunning Pseudo Code

---

**Algorithm 1:** HTML Pruner

---

**Data:** tree $tree$, kept elements $kept$, recursion count $rcc$, max depth $d$, max children $mc$, max sibling $ms$

**Result:** pruned tree $tree$

1   $nodes \leftarrow [\,]$
2   **for** $t \leftarrow 0$ **to** $rcc$ **do**
3     **for** $id$ **in** $kept$ **do**
4       $node \leftarrow tree.element$ with $id$
5       $nodes$.append($node$)
6       $nodes$.append(getAnscendants($node, d$))
7       $nodes$.append(getDescendants($node, d, mc$))
8       $nodes$.append(getSiblings($node, ms$))
9     **end for**
10    $d, mc, ms \leftarrow$ update($d,mc,ms$)// make them smaller
11 **end for**
12 **for** $node$ **in** reversed($tree$) **do**
13    **if not** $node$ **in** $nodes$ **or not** ($node$ has text or attrib **or** len($node.children$) $> 1$ **or** $node$ is root) **then**
14      $tree$.remove($element$)
15    **end if**
16 **end for**

---

# B Implementation Details of AUTOWEBGLM

During the SFT phase, we set the learning rate to 1e-5 with a batch size of 32. In the DPO stage, we sample the complex task dataset 20 times. After the filtering process, we build a contractional dataset of approximately 13k. We set the learning rate for the DPO to 1e-6, the batch size to 64, and the $\beta$ parameter to 0.15. We add the SFT loss, weighted by a factor of 0.8. During the RFT stage, we collect samples from two diverse environments, MiniWoB++ and WebArena, resulting in successful datasets of approximately 66k and 2k, respectively, which underwent finetuning. The learning rate set for this stage was 1e-5, and the batch size was 32.

# C Full Results of MiniWob++

Table 6 is the per-task average success rate on 56 tasks from MiniWoB++.

# D Input Prompt

Below is our input prompt for AUTOWEBGLM training and inference:

```
<html> {html_content} </html>

You are a helpful assistant that can assist with web navigation tasks.
You are given a simplified html webpage and a task description.
Your goal is to complete the task. You can use the provided functions
    below to interact with the current webpage.

#Provided functions:
def click(element_id: str) -> None:
    """
    Click on the element with the specified id.

    Args:
        element_id: The id of the element.
```

Table 6: PER-TASK PERFORMANCE ON MINIWOB++

| Task | AUTOWEBGLM | HTML-T5-XL | WebN-T5-XL | GPT-4 | GPT-3.5-Turbo |
|---|---|---|---|---|---|
| book-flight | 0.50 | 0.99 | 0.48 | 0.00 | 0.00 |
| choose-date | 1.00 | 0.16 | 0.08 | 0.00 | 0.00 |
| choose-date-easy | 1.00 | 1.00 | 1.00 | 0.00 | 0.00 |
| choose-date-medium | 1.00 | 0.56 | 0.07 | 0.00 | 0.00 |
| choose-list | 0.15 | 0.22 | 0.16 | 0.00 | 0.00 |
| click-button | 1.00 | 1.00 | 1.00 | 0.67 | 1.00 |
| click-button-sequence | 1.00 | 1.00 | 1.00 | 0.33 | 0.00 |
| click-checkboxes | 1.00 | 1.00 | 0.22 | 0.33 | 0.00 |
| click-checkboxes-large | 0.83 | 0.90 | 0.54 | 0.00 | 0.00 |
| click-checkboxes-soft | 0.37 | 0.99 | 0.08 | 0.00 | 0.00 |
| click-checkboxes-transfer | 1.00 | 1.00 | 0.63 | 1.00 | 0.00 |
| click-collapsible | 1.00 | 1.00 | 0.26 | 0.00 | 0.00 |
| click-collapsible-2 | 0.76 | 0.93 | 0.27 | 0.00 | 0.00 |
| click-color | 0.74 | 1.00 | 0.34 | 0.67 | 0.00 |
| click-dialog | 1.00 | 1.00 | 1.00 | 0.33 | 0.00 |
| click-dialog-2 | 1.00 | 0.74 | 1.00 | 0.67 | 0.67 |
| click-link | 1.00 | 1.00 | 0.99 | 0.33 | 0.33 |
| click-menu | 1.00 | 0.37 | 0.41 | 0.00 | 0.50 |
| click-option | 1.00 | 1.00 | 0.87 | 0.67 | 0.00 |
| click-pie | 1.00 | 0.96 | 0.51 | 0.67 | 1.00 |
| click-scroll-list | 0.57 | 0.99 | 0.98 | 0.00 | 0.00 |
| click-shades | 1.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| click-shape | 0.64 | 0.79 | 0.24 | 0.00 | 0.00 |
| click-tab | 1.00 | 1.00 | 0.57 | 0.00 | 0.67 |
| click-tab-2 | 1.00 | 0.94 | 0.57 | 0.00 | 0.00 |
| click-tab-2-hard | 1.00 | 0.88 | 0.12 | 0.33 | 0.00 |
| click-test | 1.00 | 1.00 | 1.00 | 1.00 | 0.00 |
| click-test-2 | 0.93 | 1.00 | 1.00 | 1.00 | 1.00 |
| click-widget | 1.00 | 1.00 | 1.00 | 1.00 | 0.00 |
| count-shape | 0.65 | 0.67 | 0.64 | 0.00 | 0.00 |
| email-inbox | 1.00 | 1.00 | 0.38 | 0.00 | 0.33 |
| email-inbox-forward-nl | 1.00 | 1.00 | 0.33 | 0.00 | 0.00 |
| email-inbox-forward-nl-turk | 1.00 | 1.00 | 0.23 | 0.00 | 0.00 |
| email-inbox-nl-turk | 1.00 | 0.99 | 0.20 | 0.67 | 0.00 |
| enter-date | 1.00 | 1.00 | 0.89 | 0.66 | 0.00 |
| enter-password | 1.00 | 1.00 | 0.72 | 0.67 | 0.00 |
| enter-text | 1.00 | 1.00 | 0.89 | 0.67 | 0.00 |
| enter-text-dynamic | 1.00 | 1.00 | 1.00 | 0.00 | 0.00 |
| enter-time | 0.00 | 1.00 | 0.00 | 0.00 | 0.00 |
| focus-text | 1.00 | 1.00 | 1.00 | 0.00 | 0.00 |
| focus-text-2 | 1.00 | 1.00 | 1.00 | 0.00 | 1.00 |
| grid-coordinate | 1.00 | 1.00 | 1.00 | 1.00 | 0.33 |
| guess-number | 1.00 | 0.13 | 0.00 | 0.00 | 0.00 |
| identify-shape | 1.00 | 1.00 | 0.88 | 0.67 | 0.00 |
| login-user | 1.00 | 1.00 | 0.82 | 0.33 | 0.00 |
| login-user-popup | 0.63 | 1.00 | 0.72 | 0.33 | 0.00 |
| multi-layouts | 1.00 | 1.00 | 0.83 | 0.33 | 0.00 |
| multi-orderings | 1.00 | 1.00 | 0.88 | 0.67 | 0.00 |
| navigate-tree | 1.00 | 0.99 | 0.91 | 0.33 | 0.00 |
| search-engine | 1.00 | 0.93 | 0.34 | 0.67 | 0.00 |
| social-media | 1.00 | 0.99 | 0.20 | 0.33 | 0.00 |
| social-media-all | 0.90 | 0.31 | 0.21 | 0.00 | 0.00 |
| social-media-some | 0.76 | 0.89 | 0.42 | 0.00 | 0.00 |
| tic-tac-toe | 0.74 | 0.57 | 0.48 | 0.00 | 0.00 |
| use-autocomplete | 0.85 | 0.97 | 0.02 | 1.00 | 0.67 |
| use-spinner | 1.00 | 0.07 | 0.07 | 0.00 | 0.00 |
| Average | **0.893** | 0.856 | 0.484 | 0.321 | 0.134 |

```python
    """

def hover(element_id: str) -> None:
    """
    Hover on the element with the specified id.

    Args:
        element_id: The id of the element.
    """

def select(element_id: str, option: str) -> None:
    """
    Select an option from a dropdown.

    Args:
        element_id: The id of the element.
        option: Value of the option to select.
    """

def type_string(element_id: str, content: str, press_enter: bool) ->
    None:
    """
    Type a string into the element with the specified id.

    Args:
        element_id: The id of the element.
        content: The string to type.
        press_enter: Whether to press enter after typing the string.
    """

def scroll_page(direction: Literal['up', 'down']) -> None:
    """
    Scroll down/up one page.

    Args:
        direction: The direction to scroll.
    """

def go(direction: Literal['forward', 'backward']) -> None:
    """
    Go forward/backward

    Args:
        direction: The direction to go to.
    """

def jump_to(url: str, new_tab: bool) -> None:
    """
    Jump to the specified url.

    Args:
        url: The url to jump to.
        new_tab: Whether to open the url in a new tab.
    """

def switch_tab(tab_index: int) -> None:
    """
    Switch to the specified tab.

    Args:
        tab_index: The index of the tab to switch to.
    """

def user_input(message: str) -> str:
    """
```

```
        Wait for user input.

        Args:
            message: The message to display to the user.

        Returns: The user input.
        """

def finish(answer: Optional[str]) -> None:
    """
    Finish the task (optionally with an answer).

    Args:
        answer: The answer to the task.
    """

#Previous commands: {previous_commands}

#Window tabs: {exist_window_tabs_with_pointer_to_current_tab}

#Current viewport (pages): {current_position} / {max_size}

#Task: {task_description}

You should output one command to interact to the currrent webpage.
You should add a brief comment to your command to explain your
    reasoning and thinking process.
```

## E    Data Construction Prompt

Data construction prompt for web recognition description:

```
I want you to act as a Website Reader. Your objective is to explain a
    website's purpose and usage, given the website's text. Your
    explanation should cover all of the website's primary functions.
    DO NOT GUESS the purpose of the website, you SHOULD output \"None
    \" If you are not PRETTY sure about the purpose of the website.
    Note that you should only answer the purpose or usage within 20
    words.

#Website Text#:
{html_content}

#Purpose#:
```

Data construction prompt for simple task task:

```
HTML:
{html_content}

I want you to act as a task generator that can help generate Task-
    Operation pairs.
Based on the above HTML webpage, I will give you a specified operation
    . Your goal is to come up with a ONE-STEP task that the specified
    operation can solve.
Your answer SHOULD be in the following format:

Task: {Generated one-step task}

Operation: {The right operation to solve the task}

Intention: {The intention and thinking in your operation}

NOTICE:
```

```
1. Your generated task should not be too SIMPLE, NAIVE
2. You can only do \#type\# on <input> and <textarea>
```

Data construction prompt for complex task trace intent:

```
User's overall task: {task_description}

User's actions: {annotated_action_trace}

Based on this information, deduce the intent behind each of the user's
    actions. Your response should be structured as follows:
Intent of Step 1: [Describe the intent of the user's first action from
    the user's first-person perspective]
Intent of Step 2: [Describe the intent of the user's second action
   from the user's first-person perspective]
... and so on.
Note: Your response should have the same number of lines as the number
    of user actions. The number of user actions in this task is {
   number_of_steps_in_action}.
```

## F Annotation Details

The annotation process was performed by 20 annotators for one month using the Google Chrome browser with our plugin installed to record their actions on assigned websites. The annotators first visited the target websites and checked whether the website descriptions matched the actual tasks. They then evaluated the tasks for clarity, relevance, achievability, complexity, and subjectivity, skipping those that didn't meet the criteria. They carefully recorded each step during a task, including any login or captcha steps. For tasks that required an answer, the annotators manually edited the responses. If a task was not doable, they could modify its description or abandon it. We provide the demonstration of our plugin for annotation in Figure 8 and the annotation documentation in Table 7.

Table 7: Annotation Documentation.

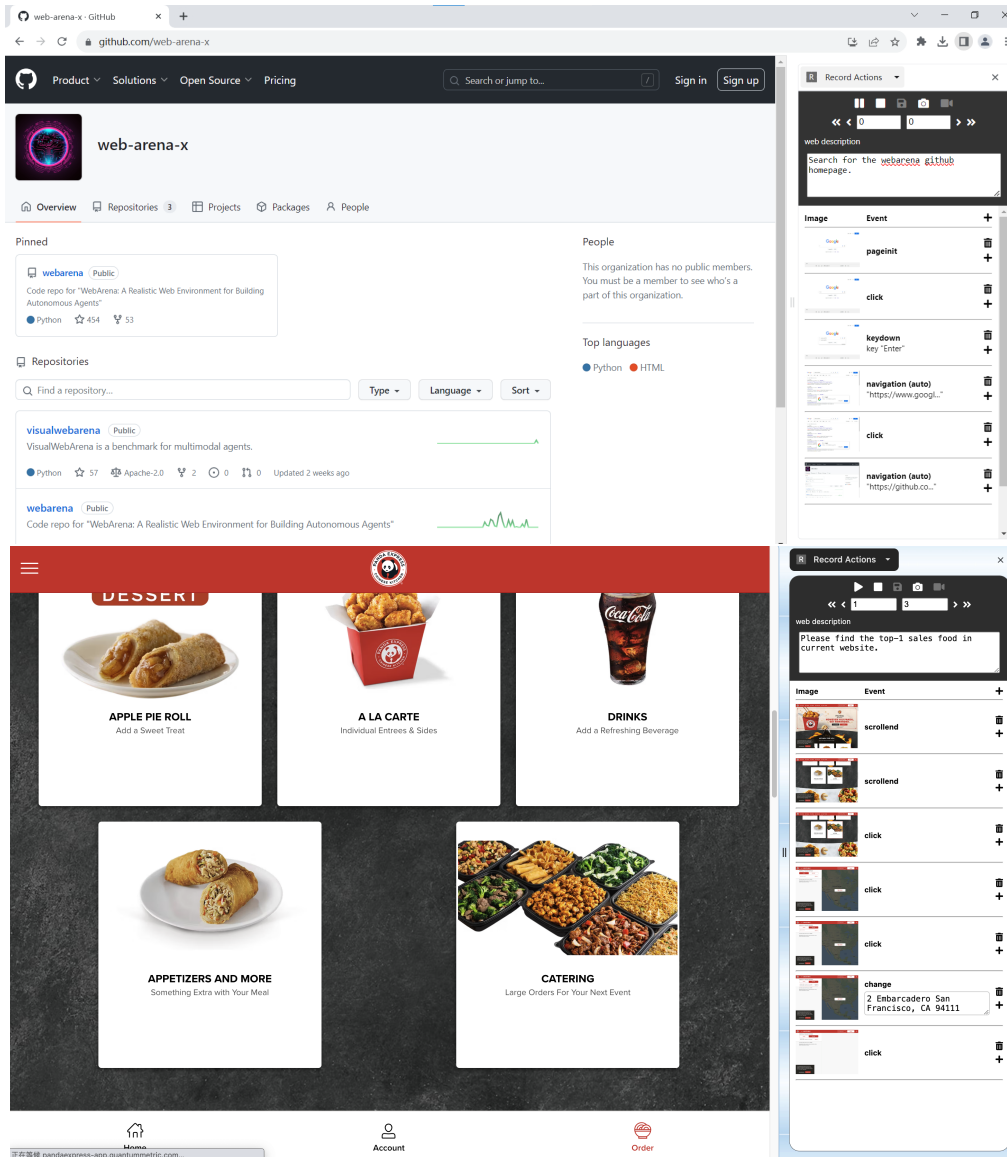| | |
|---|---|
| Annotation Target | We collected several websites and generated 50 unique tasks for each. For each website, annotators must perform operations according to the task descriptions to complete the requirements. |
| Installation of Plugin | We recommend using Google Chrome for annotation. |
| 1 | Download the plugin code zip package. |
| 2 | Download the annotation data and place the folder in the directory where your plugin code is located. |
| 3 | Start Chrome, go to More > More Tools > Extensions, enable Developer mode and load the unpacked extension by selecting the unzipped plugin code directory. |
| Plugin Introduction | To learn how to use the plugin in detail, please refer to the introduction video. The plugin has three main functions: Start/Pause Recording, End Recording, and Record Archive. |
| | Navigation between websites and tasks is made easier with the use of '«' and '»' for previous and next websites, and '<' and '>' for previous and next tasks. The plugin also provides input fields for direct navigation. |
| 1 | Click the trash bin icon to reset the environment. |
| 2 | When encountering login or robot authentication, add an authorization record by clicking the plus (+) icon one step before the corresponding operation. |
| 3 | If the task requires an answer, click the plus (+) icon in the appropriate place to insert and edit the answer. |
| 4 | If there is no response when you type content into an input field, click the plus sign and go to the appropriate location to insert and edit the input content (note that the input field is the last object to be clicked, so make sure the last click was on an input field). |
| 5 | By default, the navigation record is marked as automatic. For manual navigation, such as typing a URL or going to the previous or next page, click the word "Navigation" in this record to change it to a manual navigation record. |
| 6 | Navigation operations are not recorded until the page is fully loaded. Wait for the navigation record to appear before you continue. |
| 7 | If you need an element on the page to hover, you can click on that element and change the "click" event to a "mouseover" record. |
| Annotation Process | Here is a checklist for your reference. |
| 1 | Enter the assigned website number in the designated input box and navigate to the corresponding website. |
| 2 | First, verify if the auto-generated website description aligns with the website's functionality. The description need not be detailed, but it should be generally correct. If it matches, proceed with annotation; otherwise, skip the website. |
| 3 | Task descriptions are evaluated to determine their feasibility based on criteria such as the precision of the description, whether the task is typical of the average user of the site, and its overall feasibility. Tasks that do not involve complex login or registration operations and those that do not have too many subjective conditions are preferred. If the task has minor issues, such as an unachievable condition or involving subjective conditions, you may modify it directly in the task editing box. |
| 4 | Start executing operations according to the task description, paying attention to whether there are missing or redundant operations, and handle login or robot verification as specified. After completing the task, click "End" and then "Record Archive" to save. |
| 5 | Once annotation is complete, please upload the contents of the tasks directory. |
| Additional Annotation | Design ten tasks for each website, with each task containing at least 15 effective steps. An operation is considered effective only if it is a click, mouseover, or input change that contributes to achieving the goal. Remove any erroneous operations from the record. |
| 1 | Enter the assigned website number in the website input box and navigate to that website. |
| 2 | Briefly browse the website to determine possible tasks. |
| 3 | Choose a task number. The task box will not contain any data at this point. |
| 4 | Start recording operations. |
| 5 | After completing the operations, fill in the completed task in the task box. |
| 6 | Save the result. |

Figure 8: Demonstration of annotation plugin.

# G  Demonstration

## G.1  Weather Report

The targeted task to be executed is "What is the weather like today?". The actual execution steps can be summarized as follows:

- Step1: Type SearchBar "todays weather report"
- Step2: Click SearchButton
- Step3: Click SearchBar
- Step4: Click DateButton
- Step5: Answer

As Figure 9 shows, we end up on the webpage with a local weather report. We obtained detailed information about today's weather as the answer, effectively completing the target task.

## G.2  Shopping Advice

The targeted task to be executed is "Help me pick a Christmas gift for kids". The actual execution steps can be summarized as follows:

- Step1: Type SearchBar "Christmas gift for kids"
- Step2: Click SearchButton
- Step3: Click "All" tag in the category selection
- Step4: Click SearchBar
- Step5: Click The first product on the result page
- Step6: Answer

As Figure 10 shows, we ultimately landed on a camera product page, where we obtained a recommendation for that camera as the answer, essentially completing the task.

## G.3  Searching Article

The targeted task to be executed is "Find an article about large language model". The actual execution steps can be summarized as follows:

- Step1: Type SearchBar "large language model"
- Step2: Click SearchButton
- Step3: Click Wiki
- Step4: Scroll down
- Step5: Go backward
- Step6: Scroll down
- Step7: Click The link to an article
- Step8: Answer

As Figure 11 shows, we ultimately arrived at a page featuring an article and obtained "Found a relevant article" as the answer, essentially fulfilling the task.

## G.4  Searching Tools

The targeted task to be executed is "Find a tool to solve the differential equation". The actual execution steps can be summarized as follows:

- Step1: Type SearchBar "tools to solve the differential equation"

- Step2: Click SearchButton
- Step3: Scroll down
- Step4: Click The link to an online tool
- Step5: Click ODE calculator
- Step6: Answer

As Figure 12 shows, we ultimately arrived at a page for an ODE (Ordinary Differential Equation) calculator and obtained "Found a relevant tool" as the answer, essentially completing the task.

### G.5 Knowledge Query

The targeted task to be executed is "Search and tell me some basic info about the dark matter". The actual execution steps can be summarized as follows:

- Step1: Type SearchBar "dark matter"
- Step2: Click SearchButton
- Step3: Click Wiki
- Step4: Scroll down
- Step5: Answer

As Figure 13 shows, we ultimately reached a wiki page about dark matter, obtaining some basic info as the answer, and effectively completing the task.

### G.6 Finding Pictures

The targeted task to be executed is "Help find a beautiful picture of the Pacific Ocean". The actual execution steps can be summarized as follows:

- Step1: Type SearchBar "Pacific Ocean Pictures"
- Step2: Click SearchButton
- Step3: Click A picture in the search result
- Step4: Go backward
- Step5: Click Another in the search result
- Step6: Answer

As Figure 14 shows, we ultimately reached a page displaying an image of the Pacific Ocean, obtaining "Found a picture of the Pacific Ocean for you" as the answer, effectively completing the task.

### G.7 Finding Research

The targeted task to be executed is "Search and tell me a hot area in AI research". The actual execution steps can be summarized as follows:

- Step1: Type SearchBar "areas in AI research"
- Step2: Click SearchButton
- Step3: Click A Link to a page
- Step4: Scroll down
- Step5: Answer

As Figure 15 shows, we ultimately reached a page about AI research, obtaining "Natural Language Processing(NLP)" as the answer, effectively completing the task.

### G.8  Game Recommendation

The targeted task to be executed is "I want to play a PC game. Help me choose one". The actual execution steps can be summarized as follows:

- Step1: Type SearchBar "PC game recommendations"
- Step2: Click SearchButton
- Step3: Answer

As Figure 16 shows, we ultimately reached a page of searching results of games, obtaining a recommendation as the answer, effectively completing the task.

### G.9  Playing Video

The targeted task to be executed is "Search and tell me a hot area in AI research". The actual execution steps can be summarized as follows:

- Step1: Type SearchBar "funny videos"
- Step2: Click SearchButton
- Step3: Click A Link to a vedio
- Step4: Answer

As Figure 17 shows, we ultimately reached a page playing a funny video, effectively completing the task.

### G.10  Online Shopping Assistance with Pop-Up Interruption

The targeted task to be executed is "Find and select a highly rated toaster". The actual execution steps can be summarized as follows:

- Step1: Type SearchBar "best toaster 2024"
- Step2: Click SearchButton
- Step3: Click a link from the search results leading to an online shopping site (Encounter a pop-up asking to subscribe to the newsletter.)
- Step4: Scroll down, but the interaction is blocked by the pop-up
- Step5: Answer

As Figure 18 shows, we do not reach the intended product selection. The presence of an unexpected pop-up interrupts the task execution, demonstrating the system's limitation in handling unexpected graphical elements and pop-ups. This outcome underscores the need for enhanced capabilities in graphical recognition and interaction handling, ensuring smoother navigation and task completion on web pages with complex elements.

### G.11  Knowledge Query with Hallucination

The targeted task to be executed is "Tell me some basic info about NLP". The actual execution steps can be summarized as follows:

- Step1: Answer

As Figure 19 shows, this case is a classic example of the hallucination fallacy, where the system responded directly without going through the webpage, and the response came from the hallucination rather than the webpage information

## G.12 Technological Breakthrough Summary with misinterpretation

The targeted task to be executed is "Summarize a recent technological breakthrough in renewable energy". The actual execution steps can be summarized as follows:

- Step1: Type SearchBar "latest technological breakthrough in renewable energy 2024"
- Step2: Click SearchButton
- Step3: Click a link in search results (The system selects a link to a general overview of renewable energy technologies instead of a specific article on recent breakthroughs.)
- Step4: Scroll down
- Step5: Answer

As Figure 20 shows, we do not reach the intended outcome. Instead of summarizing a recent technological breakthrough, the system provides a generalized overview of renewable energy. This outcome highlights a misinterpretation of task context, demonstrating the system's challenge in distinguishing between general information and specific, recent developments.

## G.13 Map Query with Poor Graphical Recognition

The targeted task to be executed is "Where is Beijing relative to Shanghai according to the map". The actual execution steps can be summarized as follows:

- Step1: Type SearchBar "Beijing"
- Step2: Click SearchButton
- Step3: Answer

As Figure 21 shows, we ultimately reached a page displaying a description of Beijing and the Beijing map, obtaining "Northside" as the answer. The answer is to some extent too simple. This case illustrates two minor flaws in our system: one is that it has a slight lack of understanding of the graphical interface, and the other is that it sometimes hallucinates, and the answers it gets are not always from web information.



Figure 9: Weather Report

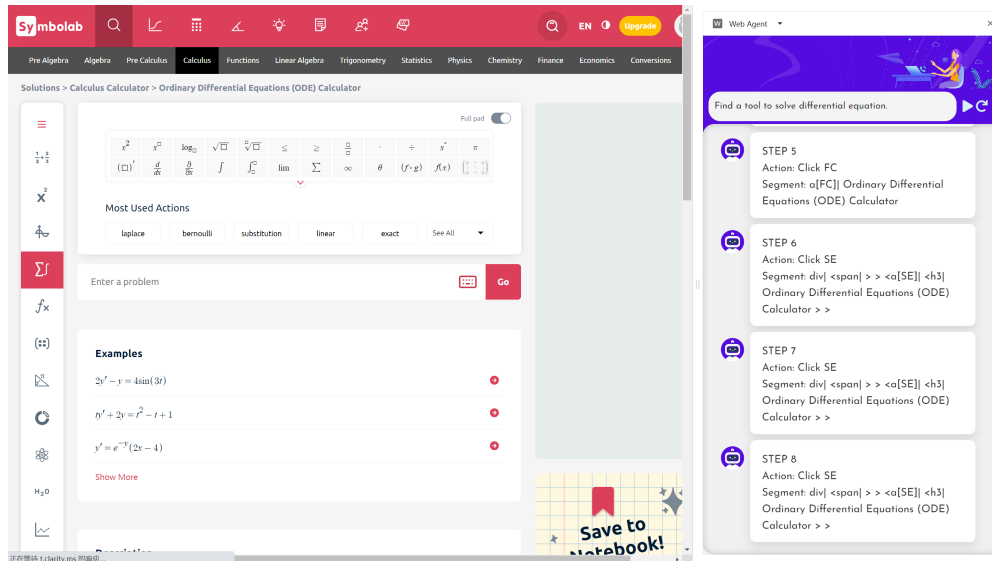Figure 10: Shopping Advice



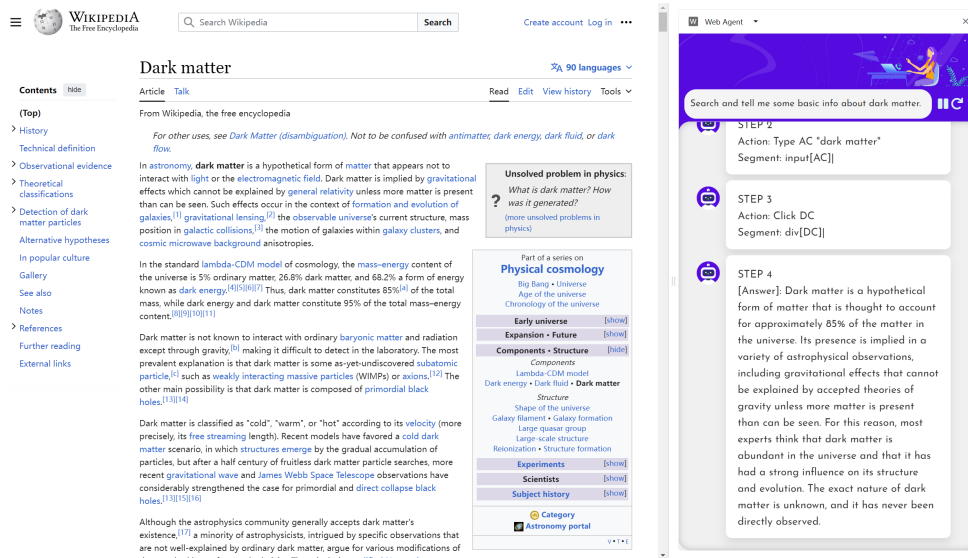Figure 11: Searching Article

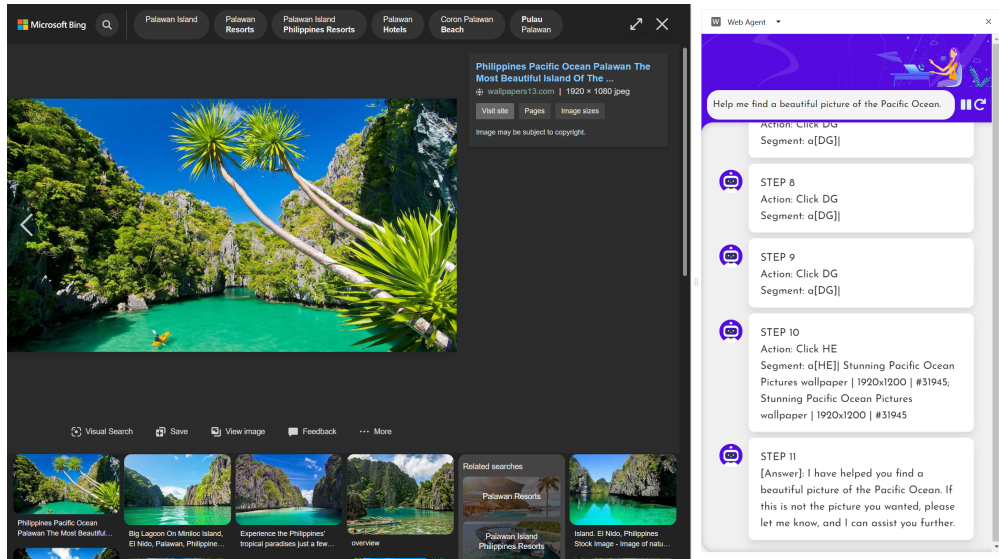Figure 12: Searching tools



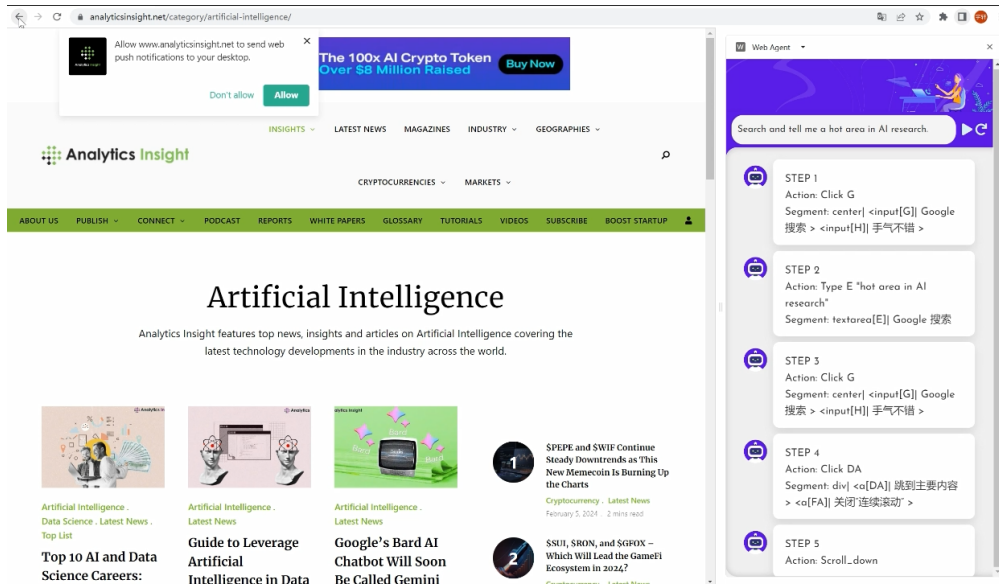Figure 13: Knowledge Query

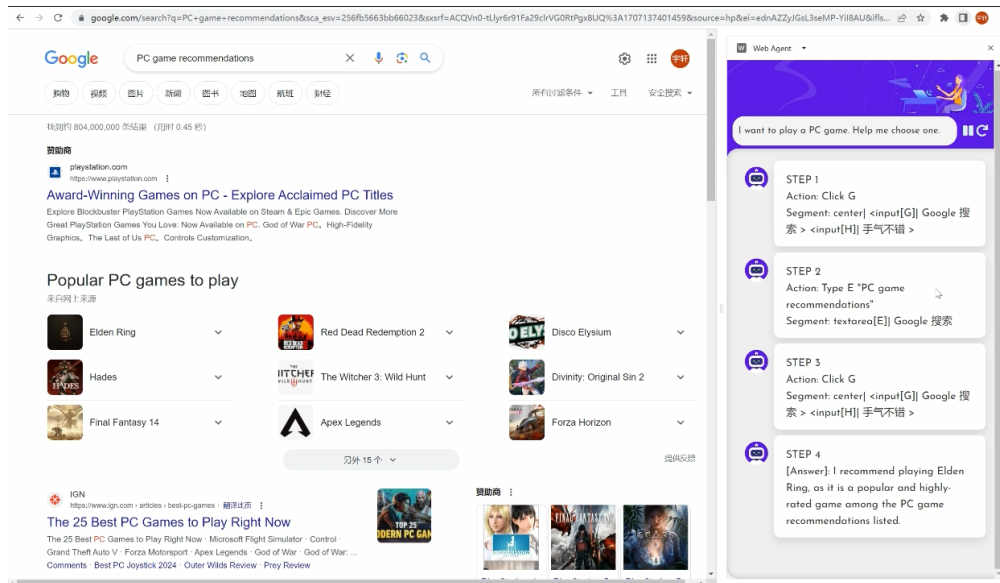Figure 14: Finding Pictures



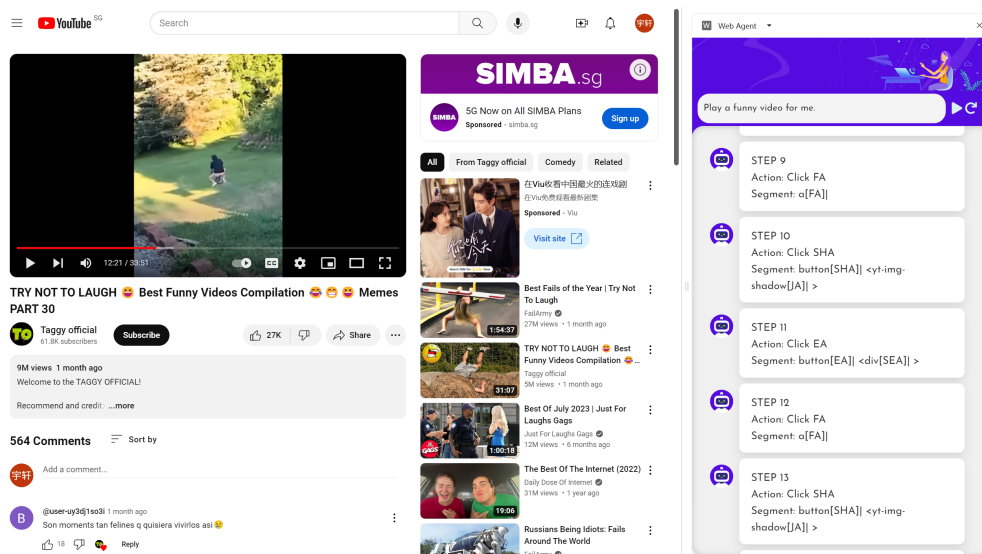Figure 15: Finding Research

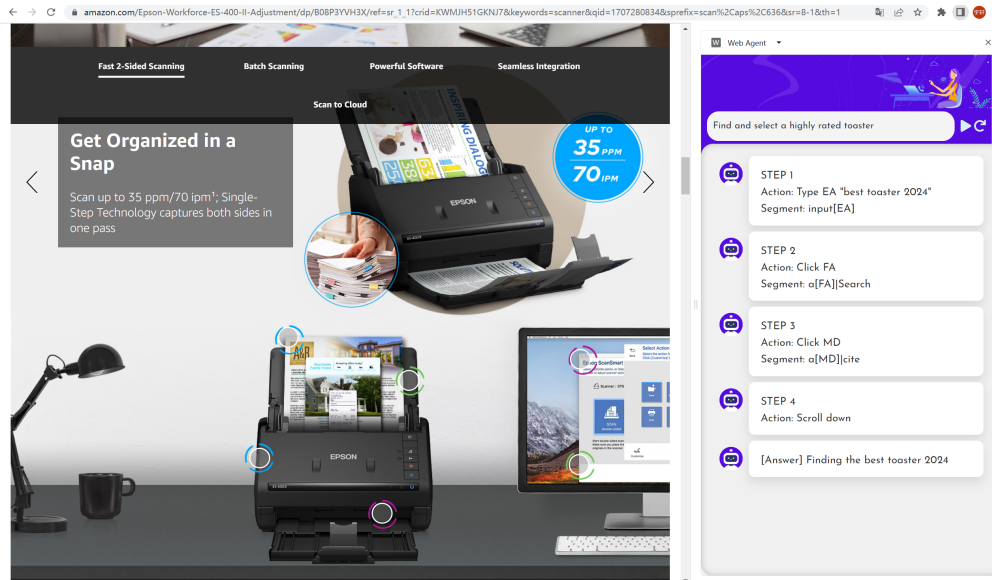Figure 16: Game Recommendation


Figure 17: Playing Video

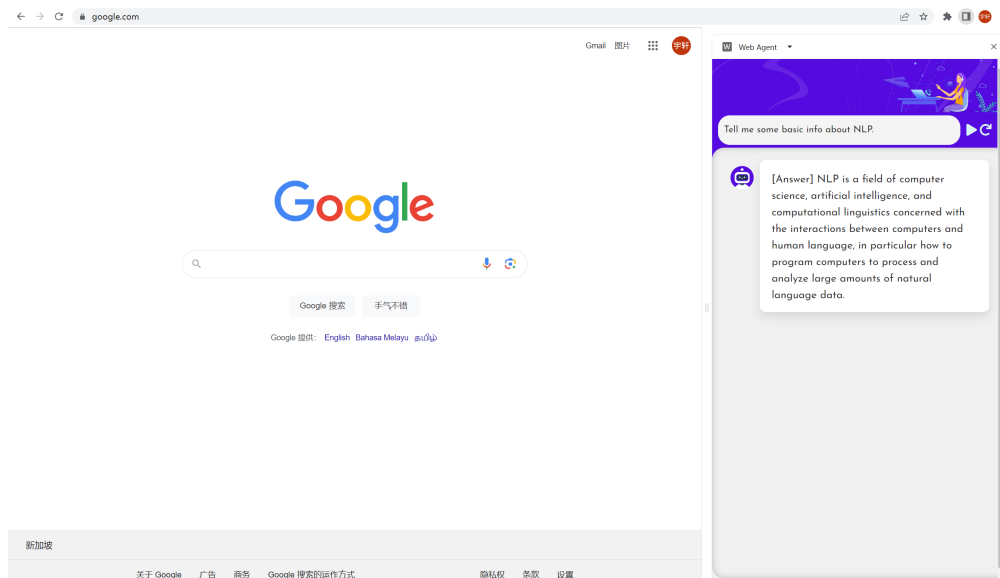Figure 18: Online Shopping Assistance with Pop-Up Interruption



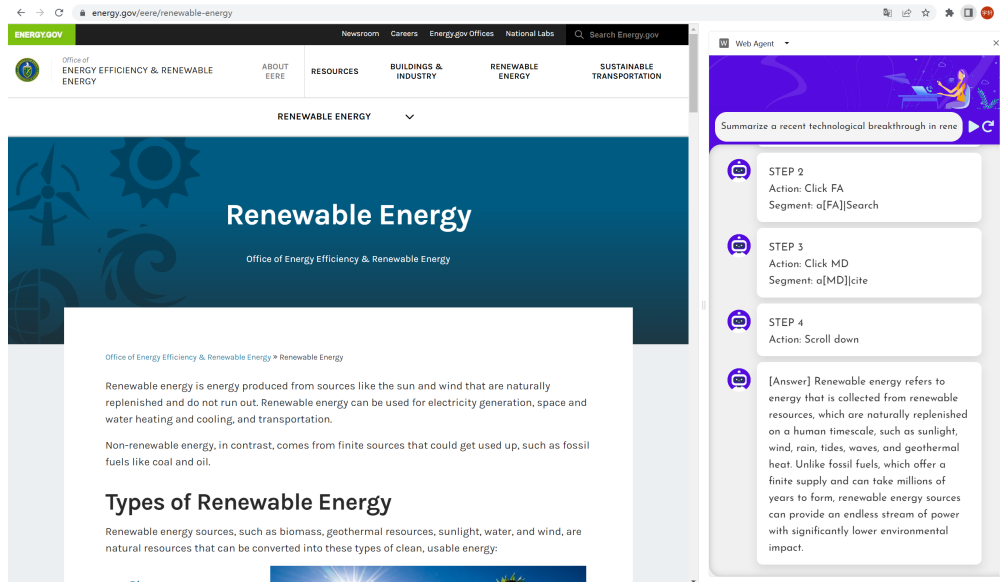Figure 19: Knowledge Query with Hallucination
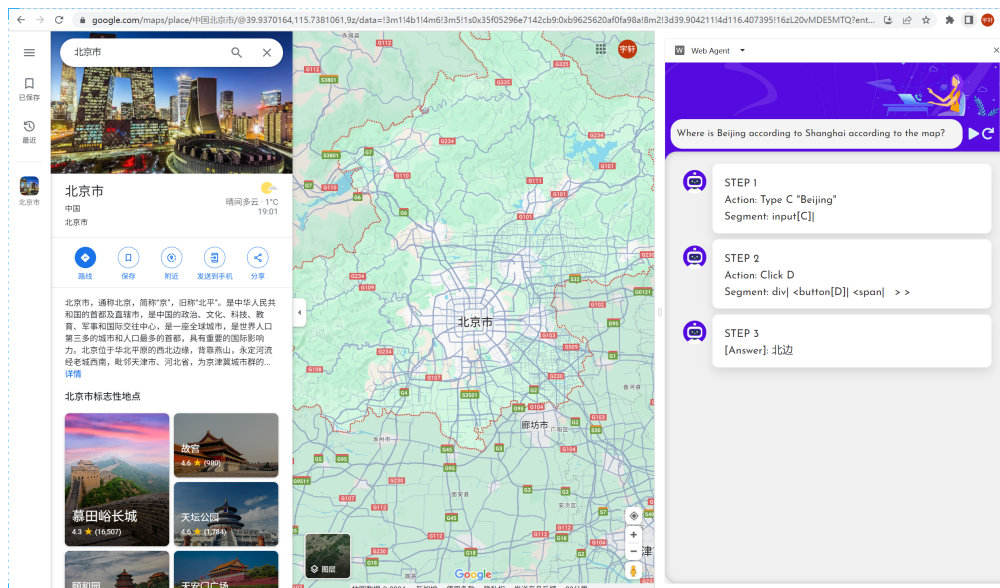
Figure 20: Technological Breakthrough Summary with misinterpretation



Figure 21: Map Query with Poor Graphical Recognition