

# NAVIQATE: Functionality-Guided Web Application Navigation

MOBINA SHAHBANDEH, University of British Columbia, Canada

PARSA ALIAN, University of British Columbia, Canada

NOOR NASHID, University of British Columbia, Canada

ALI MESBAH, University of British Columbia, Canada

End-to-end web testing is challenging due to the need to explore diverse web application functionalities. Current state-of-the-art methods, such as WebCanvas, are not designed for broad functionality exploration; they rely on specific, detailed task descriptions, limiting their adaptability in dynamic web environments. We introduce NAVIQATE, which frames web application exploration as a question-and-answer task, generating action sequences for functionalities without requiring detailed parameters. Our three-phase approach utilizes advanced large language models like GPT-4o for complex decision-making and cost-effective models, such as GPT-4o MINI, for simpler tasks. NAVIQATE focuses on functionality-guided web application navigation, integrating multi-modal inputs such as text and images to enhance contextual understanding. Evaluations on the Mind2Web-Live and Mind2Web-Live-Abstracted datasets show that NAVIQATE achieves a 44.23% success rate in user task navigation and a 38.46% success rate in functionality navigation, representing a 15% and 33% improvement over WebCanvas. These results underscore the effectiveness of our approach in advancing automated web application testing.

CCS Concepts: • **Software and its engineering** → **Software testing and debugging**.

Additional Key Words and Phrases: Large Language Models, Web Navigation, LLM Agents

## 1 INTRODUCTION

In recent years, web applications have become ubiquitous, serving as essential tools for a wide range of online activities, from e-commerce to social networking. With over 781 billion website visits globally each month [51], their popularity highlights the growing need for developers to maintain high standards of quality and functionality. Traditional manual web testing approaches, however, can be time-consuming and challenging [8], leading to the increased adoption of automated testing methodologies to streamline the quality assurance process [5, 12, 13, 19, 24, 27, 30, 44, 48, 53, 56, 64]. Despite these advances, conventional testing tools may exhibit challenges and shortcomings regarding testing coverage, potentially overlooking critical bugs and usability issues [18, 19]. The discrepancy between tests generated by conventional methods and real user interactions further compounds these challenges [63], resulting in suboptimal testing outcomes.

Web applications typically encompass a spectrum of actions, including form submissions, button clicks, and navigation through various pages. Automated testing tools for web applications encounter challenges stemming from the intricate and dynamic nature of modern web interfaces, which can feature diverse layouts, interactions, and non-deterministic states [3].

To address these challenges and mitigate the limitations of the traditional test generation methods, there has been a growing interest in leveraging deep learning (DL) [12, 13] and reinforcement learning (RL) [22, 23, 26, 27, 30, 31, 48, 64] techniques for automated testing in web applications. By assimilating insights from human testers' behavior, such automated testing approaches aim to emulate human-like interactions with web interfaces, thereby improving the comprehensiveness and effectiveness of testing. However, these DL and RL-based methodologies are not without their constraints. Challenges such as requiring a large corpus of training data, limited generalization

---

Authors' addresses: [Mobina Shahbandeh](#), University of British Columbia, Vancouver, Canada, [mobinashb@ece.ubc.ca](mailto:mobinashb@ece.ubc.ca); [Parsa Alian](#), University of British Columbia, Vancouver, Canada, [palian@ece.ubc.ca](mailto:palian@ece.ubc.ca); [Noor Nashid](#), University of British Columbia, Vancouver, Canada, [nashid@ece.ubc.ca](mailto:nashid@ece.ubc.ca); [Ali Mesbah](#), University of British Columbia, Vancouver, Canada, [amesbah@ece.ubc.ca](mailto:amesbah@ece.ubc.ca).

capabilities, and the non-deterministic nature of web applications pose hurdles to their efficacy [18].

In this context, Large Language Models (LLMs) [15–17, 40, 41, 57, 58, 62] represent a promising solution for addressing the complexities of automated testing. LLMs, such as GPT-4o, exhibit a robust understanding of both visual and textual contexts. By leveraging the capabilities of these models—particularly their ability to comprehend human knowledge and simulate human-like interactions—there is significant potential to revolutionize automated testing methodologies for web applications. LLMs have also been employed in the literature across a wide range of autonomous agents [49, 54], including web navigation agents, where they serve as the decision-making backbone [35, 44, 63? ].

However, a notable limitation of existing agents is their emphasis on testing specific, isolated tasks, rather than evaluating comprehensive functionalities. In practical applications, testing requirements frequently encompass broader functionalities such as “Search for an item” or “Add an item to the cart.” These functionalities represent more generalized actions compared to the highly specific tasks (e.g., searching for a particular item or adding a specific item to the cart) that many current agents are designed to address. This gap underscores the necessity for a more holistic approach to web application testing.

To address the aforementioned challenges, we propose a novel approach, NAVIQATE, for web application navigation, structured around three key phases: action planning, choice extraction, and decision making.

First, NAVIQATE generates concrete functionality descriptions using retrieval-augmented generation to transform abstract tasks into actionable task descriptions. From these descriptions, we extract webpage context from multi-modal sources, such as meta tags, prior actions, and screenshots, to create an abstract representation of the current state, thereby enhancing the model’s understanding of the web application’s structure and functionality. This context is leveraged to predict the next step. We then identify actionable elements on the webpage and refine them using a novel *choice ranking* system. Elements are ranked based on their semantic similarity to the predicted next step, with additional context provided by nearby visual and textual information. This ensures that only the most relevant and contextually appropriate elements are considered for action, avoiding redundant decisions and improving accuracy and efficiency. Finally, NAVIQATE selects the optimal action by combining task history, annotated screenshots, and ranked actionable elements. Visual cues help the model interpret the spatial layout and context, allowing for accurate and efficient navigation. This dynamic adaptation ensures precise task execution as NAVIQATE executes through each step of the task.

In this paper, we make the following contributions:

- A novel functionality-guided approach to web application exploration, framing web navigation as a question-and-answer task that dynamically transforms abstract functionality descriptions into concrete task steps.
- A multi-phase, multi-model methodology that enhances web navigation by integrating action planning, choice extraction, and decision making, leveraging LLMs to process textual and visual inputs for greater contextual accuracy.
- To the best of our knowledge, we are the first to enable functionality-guided navigation for automated web exploration, moving beyond predefined task-based methods.

Our empirical evaluations demonstrate that NAVIQATE attains success rates of 44.23% and 38.46% in the execution of user tasks and functionalities, respectively. Notably, these results represent a substantial improvement over the next-best baseline, WebCanvas [44], highlighting improvements of 15% and 33% in task and functionality success rates, respectively.

## 2 MOTIVATION

In end-to-end (E2E) web testing, the primary goal is often to validate the functionality of an entire system as a whole [28], rather than focusing on individual, concrete tasks. Concrete tasks require detailed parameters and predefined steps, which becomes impractical and unsustainable as applications grow in complexity. Defining these specific tasks for every possible interaction is not scalable. In contrast, focusing on functionalities allows testers to assess the overall behavior of the application without needing to define every detail. Prioritizing functionality over predefined tasks enhances the adaptability, scalability, and coverage of the testing process, leading to more reliable and robust web applications.

To test a web application, developers often define user tasks or scenarios and execute those tasks to verify intended functionalities. An example of such tasks is shown in Figure 1, where the task being performed is the following: *Find a black blazer for men with L size and add to wishlist*. The sequence of actions a user might follow to complete this task could be: 1) Select the *search* button at the top of the page, 2) Enter the value *Blazer* into the search bar, 3) Apply filters for the *Men* category, *L* size, and *Black* color, 4) Select one of the products, and 5) Click on the *Add to Wishlist* button.

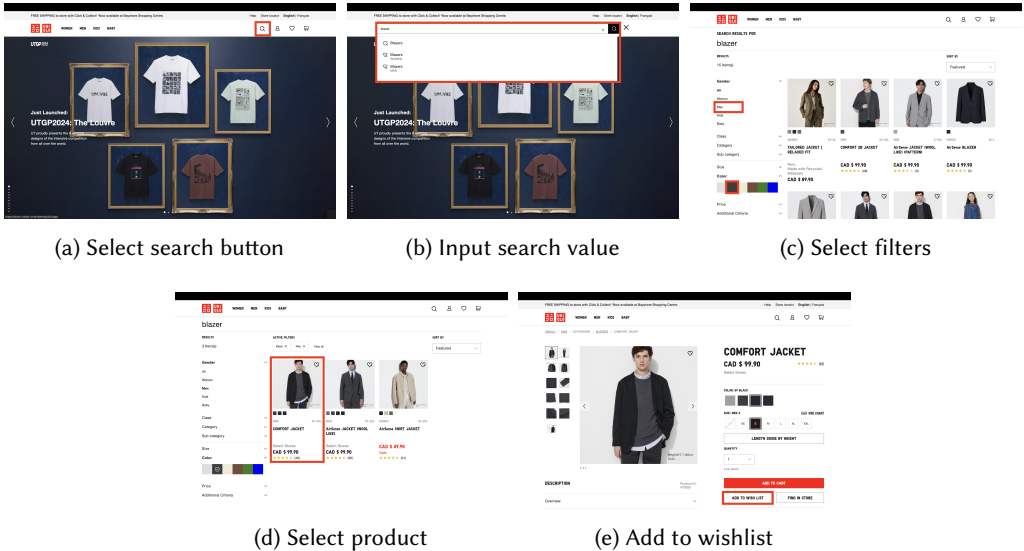


Fig. 1. Sample task execution steps for “Find a black blazer for men with L size and add to wishlist”

To clarify key terms, we define the following concepts:

**Definition 1 (User Task).** A *user task* is a detailed description of a specific operation a user can perform on a web application. This includes any specific attributes or inputs required to execute the operation. A user task consists of a finite sequence of user actions leading to the intended outcome.

Each defined task can be decomposed into two components:

- **Functionality ( $\mathcal{F}$ ):** An abstract description of the task, without specific details.
- **Task parameters ( $\mathcal{P}$ ):** A set of values required to execute the specific task, including inputs, attributes, filters, or conditions.

Each functionality  $\mathcal{F}$  represents a broad group of specific user tasks. Rather than focusing on detailed parameters, a functionality provides an abstract view of the intended action.

For example, a motivating task could consist of:

- Adding a specific item of a specific size to the wishlist as the functionality ( $\mathcal{F}$ ).
- *Blazer*, *Men*, *L*, and *Black* as the task parameters ( $\mathcal{P}$ ).

Each task defined this way can be executed through numerous variations. For instance, the task might be completed by navigating to the men’s clothing category, selecting the blazer category, applying filters, and then adding one of the resulting products to the wishlist. Alternatively, the user might search for a blazer, apply filters, and then add one of the resulting products to the wishlist.

There are several advantages to exploring functionalities, such as *Add a specific type of clothing item in a particular size to a wishlist*, over specific tasks, such as *Find a black blazer for men with L size and add to wishlist*. By understanding how to execute a functionality, the tool can gain the flexibility to perform any specific task related to that functionality, without needing to learn each task individually. This flexibility is especially valuable in applications such as virtual assistants or E2E testing, where developers only need to reference available features rather than specifying detailed *task parameters* for the agent to operate. As a result, task automation becomes more streamlined. Furthermore, because web applications frequently evolve, predefined user tasks often become obsolete over time, requiring regular updates [44]. Focusing on functionalities rather than concrete tasks offers greater efficiency, as invoking the same functionality in multiple test cases—such as login—often leads to repeated code fragments across all test scripts [28].

**Limitations of Existing Methods.** Recent research has seen increasing interest in automating web tasks, such as web question answering [10, 11]. Efforts have led to the development of natural language (NLP) models, vision models [12], and reinforcement learning models [30, 31, 64]. However, training these task-specific models is resource-intensive and time-consuming. Moreover, they often exhibit suboptimal performance when encountering novel application domains not seen during training.

In contrast, more recent approaches leverage LLMs for web task automation [5, 12, 19, 27, 44, 48, 56]. Most of these methods focus on workflows that employ LLMs to execute specific tasks. However, as previously noted, focusing on task abstractions offers significant advantages over individual tasks. To our knowledge, no prior work has explicitly pursued this direction.

Key challenges remain in developing these systems: (a) Generalization: handling variations in tasks, such as different product types or web layouts, (b) Context Awareness: the agent must understand the entire flow of actions rather than treating steps in isolation, (c) Dynamic Environments: web pages frequently change states (e.g., after applying a filter), requiring the agent to adapt accordingly, and (d) Holistic Testing: beyond one-off task execution, the system must evaluate the entire functionality (e.g., completing an entire shopping flow).

These challenges underscore the need for a more advanced, functionality-driven approach to web navigation and testing. Context awareness is crucial, as systems must make informed decisions that align with overall task objectives rather than isolated steps. In dynamic environments, where web applications and their data are constantly evolving, predefined tasks quickly become outdated and ineffective. Additionally, the focus on holistic testing is essential for evaluating entire functionalities rather than isolated tasks, ensuring a more comprehensive and robust testing strategy.

## 3 APPROACH

### 3.1 Problem Formulation

We formulate the web application navigation problem as a *sequential decision-making process* where the goal is to navigate through a series of web states to complete a given task. Each step involves



selecting an appropriate action from a predefined action space, based on the current state of the web page, the task at hand, and the predicted next step toward task completion.

NAVIQATE is tasked with generating a sequence of executable actions  $A = [A_1, A_2, \dots, A_n]$  to fulfill a given task  $U$ . At each state  $S_i$ , the model generates an action  $A_i$  by considering the current state, the next predicted step  $N_i$ , and the overall task objective  $U$ . Each action  $A_i$  consists of three components: the actionable element ( $a_i$ ), the action type ( $t_i$ ) (e.g., *click*, *type*, or *select*), and the input (if applicable) ( $i_i$ ).

$$A_i = f(S_i, U, N_i) = (a_i, t_i, i_i) \quad (1)$$

### 3.2 Actionable Elements

We focus on a variety of actionable elements, including navigation links, buttons (e.g., “Submit,” “Next”), input fields (e.g., search boxes), select elements, and any other elements that trigger actions (e.g., elements with event listeners).

The state  $S_i$  at time step  $i$  is represented by the set of all available actionable elements  $[a_0, a_1, \dots, a_n]$  in that state, and  $n$  represents the total number of actionable elements in that state, as well as all the actions performed before reaching this state  $[A_1, A_2, \dots, A_{i-1}]$ .

### 3.3 Action Space

At any given state  $S_i$ , the available actions are defined as:

- *click*( $a_i$ ): Clicking on element  $a_i$
- *type*( $a_i, i_i$ ): Typing input  $i_i$  into element  $a_i$
- *select*( $a_i, i_i$ ): Selecting the option indexed by  $i_i$  in element  $a_i$

For the *type* action,  $i_i$  represents the textual input, whereas for the *select* action,  $i_i$  refers to the index of the option to be selected. The challenge lies in choosing the right action from this space to move closer to the task goal at each step, while considering the context of the current state and the entire task trajectory.

To address the challenges outlined in previous sections, we introduce NAVIQATE, which approaches the web application navigation problem as a *sequential decision-making process*, where at each state  $S_i$ , it selects the optimal action  $a_i$  from the available action space with the goal of completing the task  $U$  efficiently by minimizing unnecessary actions and ensuring that all required steps are executed accurately. By integrating multi-modal inputs such as text and images, and combining task abstraction with dynamic action planning and selection, NAVIQATE navigates complex web applications while adapting to diverse user functionalities. The input to NAVIQATE consists of a *one-sentence task or functionality description*, such as the example provided in the motivating example, along with the name of the intended website, and the output is a *sequence of actions* to explore that task or functionality. Figure 2 provides an overview of the proposed approach. NAVIQATE employs a multi-phase, multi-model prompting methodology along with an option scoring system for web application navigation tasks, divided into three phases: **Action Planning**, **Choice Extraction**, and **Decision Making**.

### 3.4 Action Planning

The first phase of NAVIQATE, as depicted in Figure 2, is Action Planning. In this phase, we first aim to concretize the functionality descriptions (for the case of functionality exploration) to guide the LLM in generating required task parameters (defined in Definition 1), followed by a webpage context extraction for enhancing the LLM’s understanding of the high-level context of the current state, and a next step prediction, aimed at optimizing the available options for the next phase.

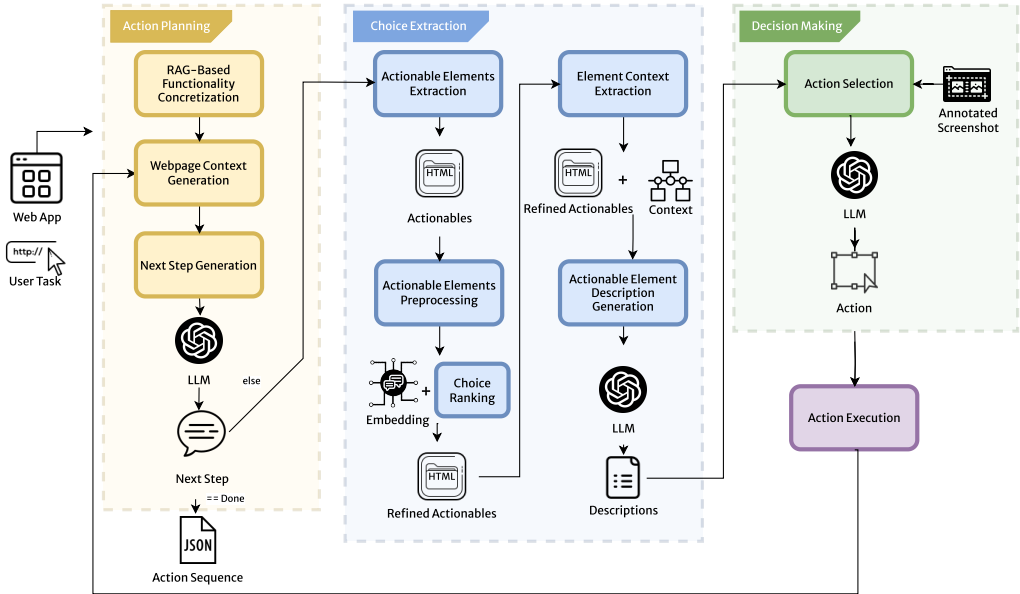


Fig. 2. Overview of NAVIQATE.

**3.4.1 Retrieval-Augmented Functionality Concretization.** In our approach to handling abstracted tasks, which involve functionality descriptions, we introduce an additional step involving the training set of a reference dataset, Mind2Web-Live [44], containing concrete task descriptions. We first generate functionality descriptions for every concrete task in the dataset, utilizing GPT-4o, to construct a reference training database, where pairs of original concrete tasks from the Mind2Web-Live training set are matched with their abstracted counterparts. For example, the task mentioned in Section 2, *Find a black blazer for men with L size and add to wishlist*, would be paired with *Add a specific type of clothing item in a particular size to a wishlist*. To generate concrete task descriptions from functionalities, we apply a Retrieval-Augmented Generation (RAG) method, which is shown to be effective in different code-related tasks [4, 42, 43]. For each given functionality description as an input to NAVIQATE, we retrieve the three most similar samples from the reference training database. The similarity measurement is based on cosine similarity between the text embeddings of the abstracted task and the tasks in the reference dataset. By converting the text of each task into high-dimensional vectors, cosine similarity allows us to measure the angular distance between these vectors, effectively capturing the semantic closeness of the tasks. This ensures that the retrieved samples are not just similar in surface-level wording but are also semantically aligned with the abstracted task. These retrieved samples are then used to prompt the language model to generate the concrete equivalent of the given abstracted task. The retrieval process ensures that no training data is directly reused during the test phase. Although the reference dataset is derived from the training data, the test dataset comprises distinct samples, separate from those used in training.

The generation of concrete tasks serves to improve the success rate of the agent’s navigation by providing clearer inputs. However, our method is designed with flexibility, allowing it to adapt to the task description as closely as possible when navigating the web application. It continuously predicts the most appropriate action at any given moment and can redirect itself as needed. Even in cases where the exact concrete version of the task is not available, the method still attempts to select the most similar possible actions to achieve the desired outcome. It is important to emphasize

that generating a concrete task from an abstract one differs from working with a fully predefined concrete task, as the former involves predicting the closest actions based on abstract descriptions, rather than relying solely on a given explicit instruction, which eliminates the need for defining detailed task descriptions.

**3.4.2 Webpage Context Generation.** To enhance the LLM’s understanding of the web application’s state at each step, we implement a webpage context generation process. This process generates a structured and abstract overview of the current webpage. To extract this context, we provide the LLM with the meta tag description of the website, the context from the previous state, the action that led to the current state, and a screenshot of the current webpage. The prompt instructs the model to describe the webpage’s context abstractly, focusing on the webpage’s primary purpose and the webpage’s sub-functionalities. An example of an LLM-generated context for the landing page of *Uniqlo* (the first state of our motivating example) is shown in [Listing 1](#). By generating this abstracted context for each state, we give the LLM a clearer understanding of the web application’s flow, thereby improving the accuracy and relevance of its responses as it navigates through the web application.

```
1 {
2   "context": "This page is the landing page for an online clothing store.",
3   "sub_functionalities": [
4     "users can browse different categories of clothing for various age groups",
5     "users can view featured products and promotions",
6     "users can navigate to other sections such as help and store locator",
7     "users can switch between different languages",
8     "users can search for specific items",
9     "users can access their account, wishlist, and shopping cart"
10  ]
11 }
```

Listing 1. Example of a generated contextual representation for a webpage

**3.4.3 Next Step Generation.** We generate a prediction of the next step,  $N_i$ , based on the task description, history of actions, and the current webpage context ([Section 3.4.2](#)). The model outputs either “Done”, indicating task completion and ending navigation, or a sentence describing the next step. In the motivating example, the initial prompt could return “Click on the ‘Clothing’ category” or “Type ‘Blazer’ in the search bar,” as tasks can have multiple valid paths. In the motivating example, the user performs the second option.

The predicted next step description serves as a guide for the LLM to select the most appropriate actionable items and inputs for the action. Additionally, incorporating the next step helps to sort and refine the list of actionables based on their semantic similarity to the generated next step. This process enhances the accuracy and relevance of the subsequent actions, ensuring that the LLM can effectively navigate and interact with the web application.

### 3.5 Choice Extraction

As shown in [Figure 2](#), the next phase, Choice Extraction, aims to prepare a list of options for the tool to choose from in the next phase. It involves extraction, preprocessing, refinement, context generation, and finally description generation of actionable elements. To extract actionable elements, we utilize the tag names and event listeners of web elements in the current state to determine if they are considered actionable as defined in [Section 3.2](#). For instance, in the motivating example discussed in [Section 2](#), the actionable elements in the first state (the landing page) include, but are not limited to, all navigation buttons, the search button, and the carousel’s navigation arrows.

```

1 {
2   "0": {
3     "outerHTML": "<a class='fr-global-nav-item px-s' href='/ca/en/men' ><div><span>men</span></div></a>",
4     "neighbours": ["MEN"]
5   },
6   "1": {
7     "outerHTML": "<input type='search' id='searchInput' autocomplete='off' class='fr-searchform-input' name='searchTerm'>",
8     "neighbours": ["Search by keyword"]
9   }, ...
10 }

```

Listing 2. Example of an actionable element’s contextual representation

**3.5.1 Actionable Elements’ Preprocessing and Ranking.** Following the extraction of the actionable elements, we undertake a ranking process for optimization purposes. This process involves calculating the semantic similarity score ( $score_{a_j}$ ) between the inner text of each actionable element ( $a_j$ ) as defined in Equation 2, or the outer HTML of the element if no inner text exists, and the textual description of the next step ( $N_i$ ), which was generated in the initial workflow step. For this task, we employ cosine similarity of the embeddings of the two texts, denoted by  $E_{a_j}$  and  $E_{N_i}$ , representing the textual embedding of the actionable element and the next step, respectively.

$$score_{a_j} = \text{cosine\_similarity}(E_{a_j}, E_{N_i}) \quad (2)$$

If an actionable element has already been selected in previous steps, we reduce its score by half. This adjustment decreases the likelihood of the element being selected repeatedly by the LLM, thereby preventing repetitive selection and potential overemphasis on the same element. After this refinement step, we sort the elements by their scores (*Choice Ranking*) and retain the top-K actionable elements, to form a list of *Refined Actionables*.

Moreover, we incorporate a preprocessing step to refine the HTML representation of actionable elements by removing certain attributes and child elements of *style*, *path*, and *svg*. Additionally, we truncate excessively large HTML elements after the preprocessing step to optimize the data we pass to the LLM.

**3.5.2 Actionable Elements’ Context Extraction.** In the next step, we extract contextual information for each element in the refined list of actionable elements and append it to its outer HTML. This involves identifying the visual neighbors of each element and incorporating their inner text. By calculating Euclidean distances within a defined threshold, we select the five closest elements to enhance the LLM’s understanding of an element’s context. For example, price or rating information near a product page link is included, providing context the link alone cannot offer. Additionally, available options for select elements are appended. An example is shown in Listing 2, where a link to the Men’s section is paired with the neighboring text “MEN,” and the search input field with “Search by keyword.”

**3.5.3 Actionable Elements’ Description Generation.** Subsequently, we prompt the LLM with the list of actionable elements and their contextual information to generate a one-sentence description for each element. We invoke the LLM to generate a description of the element’s functionality, not its appearance or style, given the textual content of its visual neighbours. We perform this step in batches of 10 elements to optimize the process. An example of an output generated by the LLM for the example actionables shown in Listing 2 is shown in Listing 3. We use a cost-effective model, GPT-4o mini, for this task, as it involves a high token count and is relatively simple for a language model to handle. For all other tasks, we use GPT-4o to take advantage of its advanced

```

1 {
2   "0": "A link to the men's section of the website.",
3   "1": "An input field with search functionality",
4   ...
5 }

```

Listing 3. Example of a generated description for an actionable element

decision-making abilities. This step lays the foundation for the next phase, Decision Making, by providing a more comprehensive list of options to choose from.

### 3.6 Decision Making

The final phase of NAVIQA<sup>TE</sup> is Decision Making. In this phase, the tool selects the best action to move closer to completing the given task or functionality. As illustrated in Figure 2, this phase includes a core component that prompts the LLM with the available actions, accompanied by visual cues, to select the next action. Once the LLM selects an action, it is grounded back into the action space for execution.

**3.6.1 Action Selection and Input Generation.** In this step, we first annotate a screenshot of the current webpage with visual cues indicating the index of each actionable element, thereby guiding the LLM regarding the visual representation and location of these elements. An example of an annotated screenshot is shown in Figure 3. We sort the list of actionable elements based on the similarity score computed using Equation 2. In the annotated screenshot example, the next step is to *click on the Blazers category link*. As a result, the most semantically similar option is the link with "Blazers" as its text, which holds the 0 index in the list. By visually mapping the elements, the LLM can better interpret their layout and context, reducing ambiguity when determining the appropriate action. The graphical cues not only help the model associate descriptions with specific locations on the page but also allow it to consider additional visual attributes, such as the shape, colour, or size of the elements. These characteristics may influence the navigation, as certain actions could depend on recognizing visual attributes of an element. Incorporating these visual details enables the LLM to make more accurate predictions for interactions that rely on both spatial and graphical information, such as selecting the correct button or entering text into the right input field.

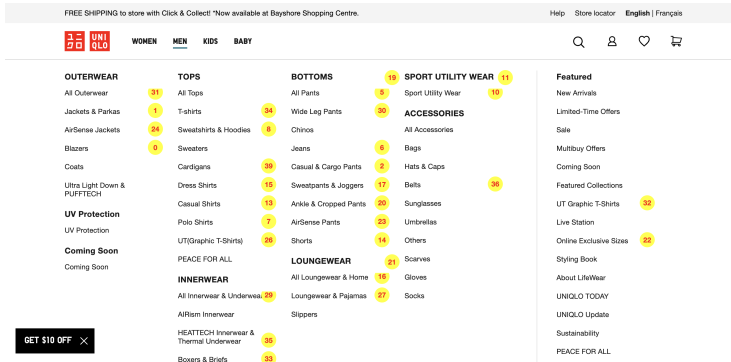


Fig. 3. Example of screenshot annotation.

The LLM is then prompted with the annotated screenshot, the list of actionable elements' descriptions and their tag names, the history of previous actions, the task description, and the next

```

1 {
2   "index": 0,
3   "action": "click"
4 }

```

Listing 4. A generated next action for *Find a black blazer for men with L size and add to wishlist*

```

1 {
2   "element": "<a class='fr-global-nav-item px-s' href='/ca/en/men' ><div><span>men</span></div></a>",
3   "xpath": "/html/body/div[1]/..",
4   "action": "click"
5 },
6 {
7   "element": "<a data-category='navi' data-action='click' data-label='/men/outerwear/blazers' href='/ca/en/men/outerwear/blazers'><div><span class='inner'>Blazers</span></div></a>",
8   "xpath": "/html/body/div[1]/..",
9   "action": "click"
10 },
11 {
12   "element": "<a role='button' data-test='filter-by-colour'><span><span class='fr-accordion-head-text'><div>Color</div></span><span class='fr-accordion-head-arrow'><span></span></span>",
13   "xpath": "/html/body/div[1]/..",
14   "action": "click"
15 } ...
16 }

```

Listing 5. Generated action sequence for *Find a black blazer for men with L size and add to wishlist* (truncated)

step description to generate the next action to be executed on the web application. Based on this input, the LLM generates the next action to be performed on the web application. It selects from the given options and determines the type of action (as explained in Section 3.3) and required arguments. NAVIQATE takes into account the need to click to activate the input field when processing *type* actions. This prevents the generation of extra click actions, making the process more efficient. An example of the LLM’s output for this step is provided in Listing 4. In this example, the LLM has selected the actionable element with an index of 0, which corresponds to a link to the “Blazers” section. This selection is then translated into executable Selenium commands, enabling interaction with the website. Finally, the selected action is executed via Selenium and appended to the action sequence for potential future reproduction.

The output of NAVIQATE is a JSON file containing a list of actions for later reproduction. Each action includes the element’s outer HTML, XPath, action type, and its arguments. An example sequence, shown in Listing 5, demonstrates NAVIQATE performing the task *Find a black blazer for men with L size and add to wishlist* by selecting links to the Men’s section, the “Blazers” section, and a filter by color.

### 3.7 Stopping Criteria

We define the task navigation stopping criteria as follows:

- The output of the Next Step Generation step is “Done”, indicating task completion as explained in Section 3.4.3, or
- There are no available actions in the current state, or
- The trajectory length exceeds a predefined limit.



To avoid potential infinite loops and to control for scenarios where the agent may get stuck in repetitive or unproductive actions during navigation, we impose a limit of 20 steps per sequence. This step limit was chosen based on empirical analysis, as most tasks can be completed well within this range since the average reference task length in the *Mind2Web-Live* dataset is 7.9.

### 3.8 Prompting

We use zero-shot prompting to prevent overfitting and ensure the model generates solutions without relying on previously seen tasks, promoting generalization.

Our prompts are designed with a system prompt that outlines the task expectations, constraints, and required output format, providing the necessary context. The user prompt supplies specific inputs, including the task description, action history, and visual context. This structure ensures the LLM fully understands the task and has the relevant details to generate accurate, context-aware actions.

## 4 EVALUATION

To assess the effectiveness of NAVIQATE, we address the following research questions:

- **RQ1:** How successful is NAVIQATE in performing a web application navigation task?
- **RQ2:** How efficient is NAVIQATE in performing a web application navigation task?
- **RQ3:** What role do different components play in the overall performance of NAVIQATE?

Domain	Subdomain Name	Sample Count
Entertainment	Event	5
	Game	3
	Movie	9
	Music	5
	Sports	9
Shopping	Auto	7
	Department	6
	Digital	6
	Fashion	3
	General	3
	Speciality	13
Travel	Airlines	5
	Car rental	1
	Ground	9
	Hotel	3
	Other	11
	Restaurant	6
<b>Total</b>		<b>104</b>

Table 1. Distribution of samples by domain and subdomain in the *Mind2Web-Live* test set [44].

### 4.1 Dataset Construction

The importance of robust datasets for training and evaluating web agents and tools has spurred research into creating reliable resources. Notable datasets in the domain of web tasks include MiniWoB++ [31], WebShop [56], RUSS [53], and Mind2Web [13]. In particular, Mind2Web offers a rich collection of real-world user interactions with web applications, comprising task descriptions and corresponding action sequences. Its subset, *Mind2Web-Live* [44], further refines this resource by excluding time-sensitive or obsolete tasks.

We utilize the test split of the *Mind2Web-Live* dataset, which comprises 104 samples, each outlining a specific user-defined task to be executed on a designated website. The websites are in three different domains, and 17 different subdomains, as shown in Table 1, with *Movie* having the

highest number of samples in *Entertainment*, *Speciality* having the most samples in *Shopping* at 13, and *Other* having the highest count at 11 in *Travel*.

To generalize each task, we abstract the descriptions by omitting task-specific details and user input information. This abstraction process is designed to generate functionality descriptions, which are more streamlined and less time-consuming to compose compared to detailed task descriptions. We employ GPT-4o to remove the task-specific details from the descriptions. We refer to this refined dataset as the *Mind2Web-Live-Abstracted*.

For the motivating example discussed in [Section 2](#), which is defined as *Find a black blazer for men with L size and add to wishlist*, the LLM generated abstraction is *Add a specific type of clothing item in a particular size to a wishlist*.

## 4.2 Implementation

NAVIQATE is implemented in Python. In our implementation, Selenium [2] serves as the backbone for automated web interactions, executing actions on web applications as directed by the workflow. Selenium is a widely-used tool for automating web browsers, enabling us to programmatically control the browser to perform tasks such as clicking buttons, filling out forms, and navigating between pages. In conjunction with Selenium, we leverage LLMs to dynamically generate and refine instructions based on the current state of the web application. For all experiments, we set the temperature parameter to 0 to ensure deterministic and well-defined answers from the LLM. The embeddings are generated using MiniLM [50]. For replaying the action sequences, we use Playwright [39], a browser automation framework that allows interaction with web elements and generates visual execution traces.

### 4.2.1 Language Models.

- GPT-4o, 2024: The GPT-4o, released by OpenAI in 2024, is engineered to provide enhanced accuracy and efficiency. We leverage this model for the more complex tasks of action planning and decision-making.
- GPT-4o MINI, 2024: GPT-4o MINI is a cost-effective model in the GPT-4o series, developed by OpenAI. We leverage this model in the simpler task of generating descriptions for web elements.

4.2.2 *LLM-only Baseline (GPT-4o)*. We implement a baseline model, utilizing GPT-4o for decision making. We only prompt the model once at each step to select an action. The model is provided with the HTML representation of actionable elements, along with the task description, to guide its action selection.

## 4.3 Evaluation metrics

**Success Rate (%)**. In the context of our study, the success rate is defined as the proportion of user tasks that the agent successfully completes on a website. Specifically, a task is considered successfully completed if the agent achieves the intended outcome. The success rate is calculated by dividing the number of successfully completed tasks by the total number of tasks assigned to the agent and then multiplying by 100 to express it as a percentage. This metric has been used in the literature for assessing the performance of web agents [44, 65].

**Trajectory Optimization Score (TOS)**. Trajectory length is defined as the number of steps or actions required to complete a given task, representing the length of the generated trajectory. Each step within the trajectory corresponds to an individual action executed by the agent. We introduce the Trajectory Optimization Score as a metric for evaluating an agent’s performance in relation to the optimal solution. The score is calculated as the ratio of the ground truth trajectory length to the

agent’s generated trajectory length for each task in the dataset. A value of 1 denotes that the agent’s trajectory aligns perfectly with the optimal path, while scores below 1 indicate inefficiencies or deviations from the ideal trajectory. To provide a comprehensive assessment, the average trajectory optimization score is computed across all tasks. For tasks in which the agent fails, a score of 0 is assigned, as higher scores correspond to better performance. This metric underscores the agent’s ability to optimize task completion.

#### 4.4 Human Evaluation

For replaying the action sequences (trajectories), we utilize Playwright [39]. Specifically, we generate a Playwright test script for each action sequence to capture execution traces in an interactive format supported by Playwright. These traces are subsequently employed for human evaluation. The authors review these traces, answering either “Yes” or “No” to assess whether the task was executed correctly. Each human evaluator (author) receives a set of execution traces, accompanied by step-by-step screenshots and a JSON file detailing the selected element (including its HTML representation and XPath), the corresponding action, and potential input to the action selected by NAVIQATE. Following a provided rubric, each author independently examines the action sequence and determines whether the task was successfully completed. The results will be considered correct only if all authors reach a positive agreement.

The Mind2Web-Live dataset includes human annotations [44], which capture the sequence of human-performed actions for each task. In their paper, WebCanvas is evaluated against these human annotations. However, as previously noted in Section 2, multiple valid approaches may exist for completing a task. Moreover, behavior cloning from expert demonstrations is prone to fragility and does not scale well [20]. In addition, annotations for functionalities are not available, and having a reference annotation for them to calculate the success rate is impractical because each functionality can correspond to multiple concrete tasks rather than a single specific task. We use the reference task length of the corresponding concrete task from the Mind2Web-Live dataset as a proxy for calculating the trajectory optimization score, with the reference task length representing one potential task aligned with the corresponding functionality.

To address this, we avoid evaluating our method against a single user-selected path. Instead, we follow a clearly defined rubric for manual evaluation. According to our rubric, a task is deemed successful if the intended functionality is explored and the correct task parameters are chosen. If additional steps are included in the sequence that do not alter the intended functionality—such as navigating to the website’s homepage after completing the task without leaving the website—we still classify the task as successful. However, these additional steps contribute to the overall trajectory length, negatively impacting the trajectory optimization score.

#### 4.5 Effectiveness and Efficiency of NAVIQATE (RQ1 & RQ2)

We evaluate our tool on both datasets: the original Mind2Web-Live dataset [44] and the abstracted dataset discussed in Section 4.1. We run the evaluation twice and report the average metrics of these two rounds. We report the metrics for both concrete and abstracted tasks independently.

*4.5.1 Comparison with State-of-the-Art Tools.* We compare our work with the WebCanvas web navigation agent introduced in [44], as this is the only existing work to our knowledge that is evaluated on the Mind2Web-Live dataset. It is worth to mention that we use the updated WebCanvas agent, which utilizes GPT-4o MINI.

Given that some tasks in the original Mind2Web dataset [13] have become unexecutable, we focus our evaluation on the refined dataset, Mind2Web-Live. To the best of our knowledge, no existing approaches specifically aim to address abstract functionalities rather than concrete tasks,

such as those found in the *Mind2Web-Live-Abstracted* dataset. Furthermore, no comparable datasets are currently available for such evaluation.

Table 2. Web task navigation performance on the Mind2Web-Live and Mind2Web-Live-Abstracted datasets.

Method	Mind2Web-Live [44]		Mind2Web-Live-Abstracted	
	SR (%)	TOS	SR (%)	TOS
WebCanvas [44]	38.46	0.39	28.84	0.16
GPT-4o	7.69	0.02	7.69	0.02
NAVIQATE	<b>44.23</b>	<b>0.58</b>	<b>38.46</b>	<b>0.56</b>

For concrete tasks in the Mind2Web-Live dataset, NAVIQATE achieves the highest SR of 44.23%, significantly outperforming WebCanvas (38.46%) and GPT-4o (7.69%). This indicates that NAVIQATE can more effectively complete the user tasks defined in the dataset. In terms of TOS, NAVIQATE also outperforms WebCanvas, achieving a score of 0.58 compared to WebCanvas’ 0.39. GPT-4o performs poorly with a TOS of 0.02, indicating it struggles with task optimization.

Similarly, for functionalities defined as abstracted tasks, NAVIQATE shows a significantly better performance, with an SR of 38.46% and a TOS of 0.56, outperforming both WebCanvas (SR: 28.84%, TOS: 0.16) and GPT-4o (SR: 7.69%, TOS: 0.02). The gap in TOS between NAVIQATE and WebCanvas is especially significant on this abstracted dataset, highlighting NAVIQATE’s ability to optimize navigation paths even when the tasks are more abstracted or generalized.

The specific tasks completed successfully differ between the abstracted and concrete datasets for the GPT-4o baseline, but the overall performance results remain consistent across both. The results demonstrate that relying solely on LLMs without any planning, context extraction, or additional guidance leads to poor navigation performance. The model struggles to consistently select the correct actions, as it lacks a structured approach to understanding the broader task context and sequence of steps.

**4.5.2 Effectiveness across domains.** The evaluation results of NAVIQATE’s performance on *Mind2Web-Live* reveal notable differences in task success at both the website and subdomain levels. The Entertainment domain achieved the highest success rate at 48.4%, and the Movie subdomain reached 88.9%. Additionally, NAVIQATE achieved a 100% success rate on certain websites, demonstrating its capability to successfully interact and navigate those sites. However, some websites showed a 0% success rate, highlighting instances where NAVIQATE was unable to complete the tasks.

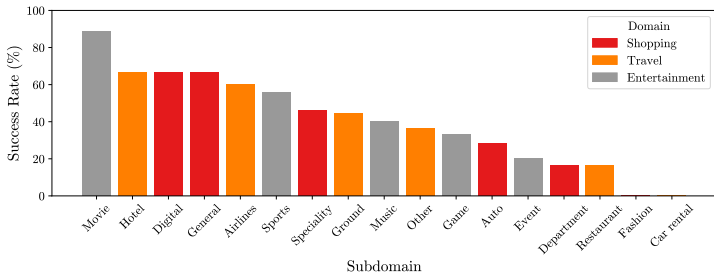


Fig. 4. Task success rate of NAVIQATE for different subdomains on the Mind2Web-Live test set.

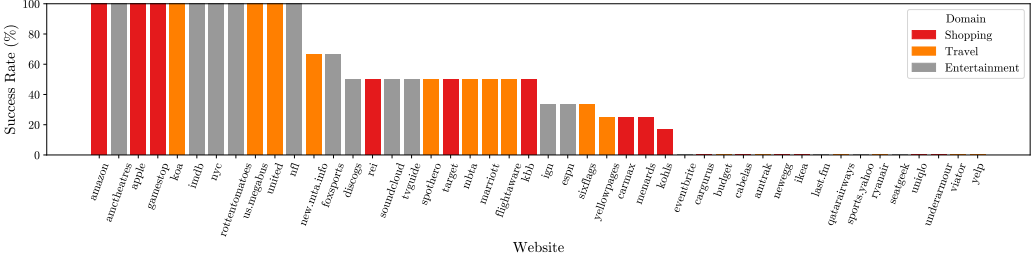


Fig. 5. Task success rate of NAVIQATE for different websites on the Mind2Web-Live test set.

#### 4.6 Ablation Study (RQ3)

**Setup.** In our ablation study, we evaluate the impact of different variations of our method to better understand the contribution of each component. Specifically, we explore three key variations:

- First, a version of our method that omits the element description generation step.
- Second, a version of our method where the agent is provided with a greater number of actionable elements (top-K) as decision-making options during action selection, compared to a version with fewer actionable elements.
- Third, a version of our method that omits the context generation and action planning steps.

By examining these variations, we aim to identify the specific components of our approach that most significantly influence the agent’s performance in navigating the web application. We select 20 random samples out of the 104 samples from the Mind2Web-Live test set for our ablation study. The rationale for selecting a smaller subset is due to the extensive number of variations in our ablation study, making it costly to run experiments across the entire dataset.

**Results.** In the first variant, we omit the element description generation (Section 3.5.3). Instead of prompting the LLM to generate descriptions of actionable elements, we directly use the textual content of these elements. We then provide a list of these textual contents, representing the actionable elements from which the agent can choose to interact in the next step of the approach. This modification leads to a decrease in the success rate, as shown in Table 3. The success rate of the original approach on the same test subset is 55%, while the success rate of the approach without element description generation decreases to 50.0%, and the TOS slightly decreases by 2%, from 0.91 to 0.89.

In the second variant, we assess the effect of the number of actionable choices (K) for the action selection. We build this variant based on the previous variant but we increase the number of actionable elements to 100 from 40. We observe no difference in the success rate but a 9% decrease in TOS. Therefore, we conclude that increasing the number of actionable choices does not improve performance. The number 40 was determined empirically and optimized for token efficiency.

In the third variant, we omit the webpage context generation (Section 3.4.2) and action planning (Section 3.4.3). This results in a decreased success rate of 45% compared to the original success rate of 55%. Additionally, the TOS declines sharply by 54%, reflecting a significant decrease in efficiency.

## 5 DISCUSSION

The results of NAVIQATE demonstrate the effectiveness of its multi-step LLM question and answering workflow in performing web application navigation tasks.

Method	SR (%)	TOS
NAVIQATE	55.0	0.91
Without element description (40 elements)	50.0	0.89
Without element description (100 elements)	50.0	0.81
Without planning	45.0	0.42

Table 3. Performance on Mind2Web-Live subset, evaluating the role of different components.

**Effectiveness of Multi-Modal Prompting.** One key finding is the effectiveness of multi-modal prompting, which allows the system to combine textual descriptions with visual inputs to improve action selection. This approach reduces ambiguities often encountered in traditional web navigation agents that rely solely on textual inputs. By incorporating annotated screenshots and global and local contexts, NAVIQATE can generate more accurate predictions of the next required action, leading to a higher task success rate compared to state-of-the-art tools like WebCanvas. Initially, we experimented with a variation of the method that excluded visual inputs (i.e., without screenshots) and observed a notable improvement when visual inputs were incorporated. This demonstrates that grounding the model in visual representations significantly enhances its predictive accuracy.

**Comparison with Baselines.** Our experimental results demonstrate that NAVIQATE achieves success rates of 44.23% and 38.46% on the Mind2Web-Live and Mind2Web-Live-Abstracted datasets, respectively. This represents a significant improvement over the WebCanvas agent, which attained success rates of 38.46% and 28.84% on the same datasets. Furthermore, NAVIQATE’s superior performance on the Mind2Web-Live-Abstracted dataset highlights its ability to effectively explore functionalities.

In many failed tasks, the WebCanvas agent would either navigate to the wrong website or move outside the correct website while attempting to complete the assigned task. In contrast, NAVIQATE remains confined within the intended website and avoids such navigation errors.

The LLM-only baseline, utilizing GPT-4o, performs notably worse, achieving success rate of only 7.69% on both the Mind2Web-Live and Mind2Web-Live-Abstracted datasets. In many failed cases, it repeatedly selects the same actionable element or chooses an incorrect action type for the element (e.g., using *type* for a link element).

Furthermore, NAVIQATE’s lower success rate in exploring functionalities compared to concrete tasks may be due to the increased difficulty of generating inputs for functionalities, as specific task parameters are not provided.

**Impact of Different Components.** The action planning and the refinement of actionable elements through context extraction and description generation was also found to contribute significantly to NAVIQATE’s performance. The ablation study revealed that omitting this step led to a reduction in task success rates, emphasizing the importance of this refinement process in ensuring accurate action generation. Interestingly, while increasing the number of actionable elements did not improve success rates, it did result in longer trajectories, indicating that having more decision options can introduce inefficiencies.

**Analysis of NAVIQATE’s Navigation Success Across Websites.** Websites with higher success rates indicate that NAVIQATE is more effective at navigating and completing tasks on these sites, likely benefiting from clearer structures or simpler interfaces. In contrast, lower success rates suggest that NAVIQATE faced challenges, such as complex layouts or technical barriers, which impeded task completion. Some websites, such as *viator*, required Captchas, preventing interaction, leading to a 0% success rate. Other task failures were due to Webdriver issues, incorrect actionable choices, or improper input generation by the model. In several cases, NAVIQATE came close to completing tasks but failed to finish them correctly, such as selecting a rating of 4 instead of 5 at



the final step of a task requiring to find “highly rated” items. Action accuracy could potentially be improved by incorporating a more extensive feedback system that detects changes in the web application in real time.

**Future Directions.** This study presents multiple opportunities for future research. One potential improvement is the implementation of an intermediate reward system, providing the model with feedback on partial task completion. This system could also be used as a new evaluation metric to improve performance measurement.

Another potential improvement could be to develop a more sophisticated feedback system, enabling the LLM to continuously monitor the effects of its actions and make more informed decisions based on web application state changes.

Moreover, a potential future direction could be the integration of a task or functionality description generation tool, such as the technique proposed by Alian et al. [6], which infers functionality descriptions for web applications. Such tools could be employed as a preliminary step to automatically generate functionality descriptions for NAVIQATE. Once generated, NAVIQATE can leverage these descriptions to perform a sequence of actions that explore the specified functionality. By incorporating this preliminary step, task descriptions could be generated automatically, leading to a more comprehensive and automated system for web application testing.

In summary, our work demonstrates the effectiveness of using LLMs for automated web application navigation and testing. Our tool, NAVIQATE, can also function as a virtual assistant for impaired users, enhancing web accessibility.

**Threats To Validity.** A potential threat to validity is the risk of LLMs attempting to mimic the examples they have seen, rather than genuinely understanding and solving the tasks. This could result in inflated performance during evaluation, as the model may simply replicate patterns observed in the training data rather than adapting to new, unseen tasks. To mitigate this, we used different datasets for RAG and evaluation, ensuring that the model was not exposed to evaluation-specific data beforehand. Additionally, we employed zero-shot prompting, which prevents the LLM from relying on previously seen samples and encourages it to generate solutions based on task descriptions alone.

A key aspect of our approach involves abstracting tasks from the Mind2Web-Live dataset. However, this abstraction process may introduce errors or inconsistencies that could impact task execution. If the abstraction is too vague, it may result in task failures. To mitigate this, we use GPT-4o to remove task-specific details from the descriptions while ensuring that the task representations remain sufficiently generalized.

The dynamic and constantly evolving nature of web applications introduces a threat to validity, as web structures, elements, and behaviors may change after the evaluation is conducted. Tasks that are executable today may become unexecutable in the future due to these changes. To address this, we rely on the Mind2Web-Live dataset, curated to reflect up-to-date web tasks.

## 6 RELATED WORK

**Web Crawling and Testing.** Web application testing [7] has advanced through automated approaches such as web crawling [37], enabling dynamic exploration and test generation [9, 38]. A foundational contribution in this area is Crawljax [37], which explores dynamic web applications by triggering actions and analyzing the DOM. Recent advancements, such as QExplore[45], apply Q-learning to improve the exploration of web applications, while Yandrapally et al.[55] address near-duplicate states detection in web applications, enhancing model accuracy and testing efficiency by eliminating redundant states. Similarly, WebEmbed [47] uses neural embeddings to abstract web page states, improving duplicate detection and coverage during web crawling. Overall, these

approaches focus on maximizing exploration by interacting with user interfaces to generate new states. In contrast, NAVIQATE adopts a functionality-guided approach to web application exploration, navigating based on a given functionality or task. This strategy allows for more structured and goal-oriented navigation, enhancing the ability to complete complex functionalities within the application.

**LLM-based Mobile and Web Testing.** Recent studies have increasingly focused on using LLMs to automate various aspects of web testing. Advances in LLMs have enabled their application in a range of testing domains, including GUI testing for both web and mobile applications [32–34]. Additionally, LLMs are being used to generate automated web form tests, simulating user interactions to validate form functionalities [5, 29].

**Web Agents.** In recent years, there has been increasing interest in developing web agents capable of determining and executing sequences of actions within web applications based on natural language instructions [22, 30, 36, 53]. The key challenge lies in creating agents that can interpret complex user instructions and translate them into appropriate web-based actions.

The advent of language models has introduced a new paradigm for web agents, removing the need for expert demonstrations [22], human-designed heuristics [64], or restricting actions to predefined sets [53]. These models leverage advanced reasoning capabilities to determine the next steps based on natural language inputs [18, 21, 46, 49]. This flexibility allows LLM-based agents to handle a wider range of web tasks, making them more versatile than previous methods.

Recent work, such as WebCanvas [44], uses LLMs to guide agents in navigating web environments. However, it relies heavily on detailed task descriptions and lacks support for broad functionalities, limiting its applications to complex tasks. In contrast, our approach supports a wide range of functionalities without requiring detailed task descriptions.

**LLMs for Software Testing.** LLMs are increasingly applied in software testing, including unit test generation [43, 52, 60], security testing [61], and defect reproduction [25]. Additionally, studies have explored their role in predicting flaky tests [14] and automating the generation and migration of test scripts [59].

## 7 CONCLUSION

In conclusion, this paper introduces NAVIQATE, a functionality-guided web application navigation tool that leverages large language models to enhance automated web testing. By framing web task navigation as a question-and-answer task and utilizing a three-phase approach with a combination of advanced and cost-efficient models, NAVIQATE significantly improves task success rates and testing efficiency. Evaluated on the Mind2Web-Live dataset, NAVIQATE demonstrated substantial improvements over the state-of-the-art, achieving a success rate of 44.23% on user tasks and 38.46% on functionalities. These results represent an improvement of 15% and 33% over the WebCanvas agent, along with a trajectory optimization score increase of 49% and 250%. NAVIQATE’s integration of multi-modal inputs and functionality-driven navigation enables it to navigate complex web tasks effectively. These contributions lay the groundwork for advancing automated web testing tools.

## 8 DATA AVAILABILITY

The implementation of our technique, NAVIQATE, along with the accompanying Mind2Web-Live-Abstracted dataset, has been made publicly available [1]. Additionally, comprehensive documentation for NAVIQATE is provided to ensure the reproducibility of our findings.

## REFERENCES

- [1] 2024. NaviQAte. <https://anonymous.4open.science/r/naviqate>.
- [2] 2024. Selenium. <https://www.selenium.dev>. Accessed: 2024-06-28.
- [3] Christoffer Quist Adamsen, Anders Møller, Rezwana Karim, Manu Sridharan, Frank Tip, and Koushik Sen. 2017. Repairing Event Race Errors by Controlling Nondeterminism. In *2017 IEEE/ACM 39th International Conference on Software Engineering (ICSE)*. 289–299.
- [4] Toufique Ahmed, Kunal Suresh Pai, Premkumar Devanbu, and Earl Barr. 2024. Automatic Semantic Augmentation of Language Model Prompts (for Code Summarization). In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*. Association for Computing Machinery, Article 220, 13 pages.
- [5] Parsa Alian, Noor Nashid, Mobina Shahbandeh, and Ali Mesbah. 2024. Semantic Constraint Inference for Web Form Test Generation. In *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis (ISSTA 2024)*. Association for Computing Machinery, 932–944.
- [6] Parsa Alian, Noor Nashid, Mobina Shahbandeh, Taha Shabani, and Ali Mesbah. 2024. A Feature-Based Approach to Generating Comprehensive End-to-End Tests. *arXiv preprint arXiv:2408.01894* (2024).
- [7] Sebastian Balsam and Deepti Mishra. 2024. Web application testing—Challenges and opportunities. *Journal of Systems and Software* 219 (2024), 112186.
- [8] Matteo Biagiola, Andrea Stocco, Filippo Ricca, and Paolo Tonella. 2019. Diversity-based web test generation. In *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE 2019)*. Association for Computing Machinery, 142–153.
- [9] Matteo Biagiola, Andrea Stocco, Filippo Ricca, and Paolo Tonella. 2020. Dependency-aware web test generation. In *2020 IEEE 13th International Conference on Software Testing, Validation and Verification (ICST)*. IEEE, 175–185.
- [10] Yingshan Chang, Guihong Cao, Mridu Narang, Jianfeng Gao, Hisami Suzuki, and Yonatan Bisk. 2022. WebQA: Multihop and Multimodal QA. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 16474–16483.
- [11] Qiaochu Chen, Aaron Lamoreaux, Xinyu Wang, Greg Durrett, Osbert Bastani, and Isil Dillig. 2021. Web question answering with neurosymbolic program synthesis. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation (PLDI 2021)*. Association for Computing Machinery, 328–343.
- [12] Qi Chen, Dileepa Pitawela, Chongyang Zhao, Gengze Zhou, Hsiang-Ting Chen, and Qi Wu. 2024. Webvln: Vision-and-language navigation on websites. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 1165–1173.
- [13] Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su. 2024. Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing Systems* 36 (2024).
- [14] Sakina Fatima, Taher A Ghaleb, and Lionel Briand. 2022. Flakify: A black-box, language model-based predictor for flaky tests. *IEEE Transactions on Software Engineering* 49, 4 (2022), 1912–1927.
- [15] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu Zheng, Ke Li, Xing Sun, et al. 2023. MME: A Comprehensive Evaluation Benchmark for Multimodal Large Language Models. *arXiv preprint arXiv:2306.13394* (2023).
- [16] Chaoyou Fu, Yuhang Dai, Yondong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. 2024. Video-MME: The First-Ever Comprehensive Evaluation Benchmark of Multi-modal LLMs in Video Analysis. *arXiv preprint arXiv:2405.21075* (2024).
- [17] Chaoyou Fu, Renrui Zhang, Haojia Lin, Zihan Wang, Timin Gao, Yongdong Luo, Yubo Huang, Zhengye Zhang, Longtian Qiu, Gaoxiang Ye, et al. 2023. A challenger to gpt-4v? early explorations of gemini in visual expertise. *arXiv preprint arXiv:2312.12436* (2023).
- [18] Hiroki Furuta, Ofir Nachum, Kuang-Huei Lee, Yutaka Matsuo, Shixiang Shane Gu, and Izzeddin Gur. 2023. Multimodal Web Navigation with Instruction-Finetuned Foundation Models. *arXiv preprint arXiv:2305.11854* (2023).
- [19] Izzeddin Gur, Hiroki Furuta, Austin V Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck, and Aleksandra Faust. 2024. A Real-World WebAgent with Planning, Long Context Understanding, and Program Synthesis. In *The Twelfth International Conference on Learning Representations*.
- [20] Izzeddin Gur, Natasha Jaques, Kevin Malta, Manoj Tiwari, Honglak Lee, and Aleksandra Faust. 2021. Adversarial Environment Generation for Learning to Navigate the Web. *arXiv preprint arXiv:2103.01991* (2021).
- [21] Izzeddin Gur, Ofir Nachum, Yingjie Miao, Mustafa Safdari, Austin Huang, Aakanksha Chowdhery, Sharan Narang, Noah Fiedel, and Aleksandra Faust. 2022. Understanding html with large language models. *arXiv preprint arXiv:2210.03945* (2022).
- [22] Izzeddin Gur, Ulrich Rueckert, Aleksandra Faust, and Dilek Hakkani-Tur. 2018. Learning to navigate the web. *arXiv preprint arXiv:1812.09195* (2018).
- [23] Miyoung Han, Pierre-Henri Wuillemin, and Pierre Senellart. 2018. Focused Crawling Through Reinforcement Learning. In *International Conference on Web Engineering*.
- [24] Javaria Imtiaz, Muhammad Zohaib Iqbal, and Muhammad Uzair Khan. 2021. An automated model-based approach to repair test suites of evolving web applications. *Journal of Systems and Software* 171 (2021), 110841.

- [25] Sungmin Kang, Juyeon Yoon, and Shin Yoo. 2023. Large language models are few-shot testers: Exploring llm-based general bug reproduction. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 2312–2323.
- [26] Andreas Kontogiannis, Dimitrios Kelesis, Vasilis Pollatos, Georgios Paliouras, and George Giannakopoulos. 2021. Tree-based focused web crawling with reinforcement learning. *arXiv preprint arXiv:2112.07620* (2021).
- [27] Hanyu Lai, Xiao Liu, Iat Long Iong, Shuntian Yao, Yuxuan Chen, Pengbo Shen, Hao Yu, Hanchen Zhang, Xiaohan Zhang, Yuxiao Dong, et al. 2024. AutoWebGLM: Bootstrap And Reinforce A Large Language Model-based Web Navigating Agent. *arXiv preprint arXiv:2404.03648* (2024).
- [28] Maurizio Leotta, Diego Clerissi, Filippo Ricca, and Paolo Tonella. 2016. Chapter Five - Approaches and Tools for Automated End-to-End Web Testing. *Advances in Computers*, Vol. 101. Elsevier, 193–237.
- [29] Tao Li, Chenhui Cui, Lei Ma, Dave Towey, Yujie Xie, and Rubing Huang. 2024. Leveraging Large Language Models for Automated Web-Form-Test Generation: An Empirical Study. *arXiv preprint arXiv:2405.09965* (2024).
- [30] Yuanchun Li and Oriana Riva. 2021. Glider: A reinforcement learning approach to extract UI scripts from websites. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1420–1430.
- [31] Evan Zheran Liu, Kelvin Guu, Panupong Pasupat, Tianlin Shi, and Percy Liang. 2018. Reinforcement learning on web interfaces using workflow-guided exploration. *arXiv preprint arXiv:1802.08802* (2018).
- [32] Zhe Liu, Chunyang Chen, Junjie Wang, Xing Che, Yuekai Huang, Jun Hu, and Qing Wang. 2023. Fill in the Blank: Context-aware Automated Text Input Generation for Mobile GUI Testing. In *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. 1355–1367.
- [33] Zhe Liu, Chunyang Chen, Junjie Wang, Mengzhuo Chen, Boyu Wu, Xing Che, Dandan Wang, and Qing Wang. 2024. Make llm a testing expert: Bringing human-like interaction to mobile gui testing via functionality-aware decisions. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*. 1–13.
- [34] Zhe Liu, Chunyang Chen, Junjie Wang, Mengzhuo Chen, Boyu Wu, Zhilin Tian, Yuekai Huang, Jun Hu, and Qing Wang. 2024. Testing the limits: Unusual text inputs generation for mobile app crash detection with large language model. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*. 1–12.
- [35] Kaixin Ma, Hongming Zhang, Hongwei Wang, Xiaoman Pan, Wenhao Yu, and Dong Yu. 2024. LASER: LLM Agent with State-Space Exploration for Web Navigation. *arXiv preprint arXiv:2309.08172* (2024).
- [36] Sahisnu Mazumder and Oriana Riva. 2021. FLIN: A Flexible Natural Language Interface for Web Navigation, In Annual Conference of the North American Chapter of the Association for Computational Linguistics (NACCL 2021). *arXiv preprint arXiv:2010.12844*.
- [37] Ali Mesbah, Arie van Deursen, and Stefan Lenselink. 2012. Crawling Ajax-Based Web Applications through Dynamic Analysis of User Interface State Changes. *ACM Trans. Web* 6, 1, Article 3 (mar 2012), 30 pages.
- [38] Ali Mesbah, Arie van Deursen, and Danny Roest. 2012. Invariant-Based Automatic Testing of Modern Web Applications. *IEEE Transactions on Software Engineering* 38, 1 (2012), 35–53.
- [39] Microsoft. 2020. *Playwright*. <https://playwright.dev> A framework for Web Testing and Automation..
- [40] Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. 2024. Large language models: A survey. *arXiv preprint arXiv:2402.06196* (2024).
- [41] Chancharik Mitra, Brandon Huang, Trevor Darrell, and Roei Herzig. 2024. Compositional chain-of-thought prompting for large multimodal models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 14420–14431.
- [42] Noor Nashid, Taha Shabani, Parsa Alian, and Ali Mesbah. 2024. Contextual API Completion for Unseen Repositories Using LLMs. *arXiv preprint arXiv:2405.04600* (2024).
- [43] Noor Nashid, Mifta Sintaha, and Ali Mesbah. 2023. Retrieval-Based Prompt Selection for Code-Related Few-Shot Learning. In *Proceedings of the 45th International Conference on Software Engineering (ICSE '23)*. IEEE Press, 2450–2462.
- [44] Yichen Pan, Dehan Kong, Sida Zhou, Cheng Cui, Yifei Leng, Bing Jiang, Hangyu Liu, Yanyi Shang, Shuyan Zhou, Tongshuang Wu, et al. 2024. WebCanvas: Benchmarking Web Agents in Online Environments. *arXiv preprint arXiv:2406.12373* (2024).
- [45] Salman Sherin, Asmar Muqheet, Muhammad Uzair Khan, and Muhammad Zohaib Iqbal. 2023. QExplore: An exploration strategy for dynamic web applications using guided search. *Journal of Systems and Software* 195 (2023), 111512.
- [46] Paloma Sodhi, SRK Branavan, and Ryan McDonald. 2023. HeaP: Hierarchical Policies for Web Actions using LLMs. *arXiv preprint arXiv:2310.03720* (2023).
- [47] Andrea Stocco, Alexandra Willi, Luigi Libero Lucio Starace, Matteo Biagiola, and Paolo Tonella. 2023. Neural embeddings for web testing. *arXiv preprint arXiv:2306.07400* (2023).
- [48] Lucas-Andrei Thil, Mirela Popa, and Gerasimos Spanakis. 2024. Navigating WebAI: Training Agents to Complete Web Tasks with Large Language Models and Reinforcement Learning. In *Proceedings of the 39th ACM/SIGAPP Symposium on Applied Computing*. 866–874.

- [49] Lei Wang, Chengbang Ma, Xueyang Feng, Zeyu Zhang, Hao ran Yang, Jingsen Zhang, Zhi-Yang Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Ji rong Wen. 2023. A Survey on Large Language Model based Autonomous Agents. *arXiv preprint arXiv:2308.11432* (2023).
- [50] Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. 2020. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. *Advances in Neural Information Processing Systems* 33 (2020), 5776–5788.
- [51] Henrique S. Xavier. 2024. The Web unpacked: a quantitative analysis of global Web usage. *arXiv preprint arXiv:2404.17095* (2024).
- [52] Zhuokui Xie, Yinghao Chen, Chen Zhi, Shuiguang Deng, and Jianwei Yin. 2023. ChatUniTest: a ChatGPT-based automated unit test generation tool. *arXiv:2305.04764* (2023).
- [53] Nancy Xu, Sam Masling, Michael Du, Giovanni Campagna, Larry Heck, James Landay, and Monica S Lam. 2021. Grounding open-domain instructions to automate web support tasks. *arXiv preprint arXiv:2103.16057* (2021).
- [54] An Yan, Zhengyuan Yang, Wanrong Zhu, Kevin Lin, Linjie Li, Jianfeng Wang, Jianwei Yang, Yiwu Zhong, Julian McAuley, Jianfeng Gao, Zicheng Liu, and Lijuan Wang. 2023. GPT-4V in Wonderland: Large Multimodal Models for Zero-Shot Smartphone GUI Navigation. *arXiv preprint arXiv:2311.07562* (2023).
- [55] Rahulkrishna Yandrapally, Andrea Stocco, and Ali Mesbah. 2020. Near-duplicate detection in web app model inference. In *Proceedings of the ACM/IEEE 42nd international conference on software engineering*. 186–197.
- [56] Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. 2022. Webshop: Towards scalable real-world web interaction with grounded language agents. *Advances in Neural Information Processing Systems* 35 (2022), 20744–20757.
- [57] Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. 2023. A survey on multimodal large language models. *arXiv preprint arXiv:2306.13549* (2023).
- [58] Shukang Yin, Chaoyou Fu, Sirui Zhao, Tong Xu, Hao Wang, Dianbo Sui, Yunhang Shen, Ke Li, Xing Sun, and Enhong Chen. 2023. Woodpecker: Hallucination correction for multimodal large language models. *arXiv preprint arXiv:2310.16045* (2023).
- [59] Shengcheng Yu, Chunrong Fang, Yuchen Ling, Chentian Wu, and Zhenyu Chen. 2023. Llm for test script generation and migration: Challenges, capabilities, and opportunities. In *2023 IEEE 23rd International Conference on Software Quality, Reliability, and Security (QRS)*. IEEE, 206–217.
- [60] Zhiqiang Yuan, Yiling Lou, Mingwei Liu, Shiji Ding, Kaixin Wang, Yixuan Chen, and Xin Peng. 2023. No More Manual Tests? Evaluating and Improving ChatGPT for Unit Test Generation. *arXiv:2305.04207* (2023).
- [61] Ying Zhang, Wenjia Song, Zhengjie Ji, Na Meng, et al. 2023. How well does llm generate security tests? *arXiv preprint arXiv:2310.00710* (2023).
- [62] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223* (2023).
- [63] Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. 2024. Gpt-4v (ision) is a generalist web agent, if grounded. *arXiv preprint arXiv:2401.01614* (2024).
- [64] Yan Zheng, Yi Liu, Xiaofei Xie, Yepang Liu, Lei Ma, Jianye Hao, and Yang Liu. 2021. Automatic web testing using curiosity-driven reinforcement learning. In *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*. IEEE, 423–435.
- [65] Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, Uri Alon, and Graham Neubig. 2024. WebArena: A Realistic Web Environment for Building Autonomous Agents. *arXiv preprint arXiv:2307.13854* (2024).