

# Feature-Driven End-To-End Test Generation

Parsa Alian  
University of British Columbia  
Vancouver, Canada  
palian@ece.ubc.ca

Noor Nashid  
University of British Columbia  
Vancouver, Canada  
nashid@ece.ubc.ca

Mobina Shahbandeh  
University of British Columbia  
Vancouver, Canada  
mobinashb@ece.ubc.ca

Taha Shabani  
University of British Columbia  
Vancouver, Canada  
taha.shabani@ece.ubc.ca

Ali Mesbah  
University of British Columbia  
Vancouver, Canada  
amesbah@ece.ubc.ca

**Abstract**—End-to-end (E2E) testing is essential for ensuring web application quality. However, manual test creation is time-consuming, and current test generation techniques produce incoherent tests. In this paper, we present AUTOE2E, a novel approach that leverages Large Language Models (LLMs) to automate the generation of semantically meaningful feature-driven E2E test cases for web applications. AUTOE2E intelligently infers potential features within a web application and translates them into executable test scenarios. Furthermore, we address a critical gap in the research community by introducing E2EBENCH, a new benchmark for automatically assessing the feature coverage of E2E test suites. Our evaluation on E2EBENCH demonstrates that AUTOE2E achieves an average feature coverage of 79%, outperforming the best baseline by 558%, highlighting its effectiveness in generating high-quality, comprehensive test cases.

**Index Terms**—Feature Inference, End-to-End Testing, Large Language Models

## I. INTRODUCTION

End-to-End (E2E) testing assesses whether various integrated components in an application work together correctly from the user interface (UI) to the back-end, by simulating real user interactions and verifying the application’s functionality. In E2E, the application is tested as a whole, in its entirety, and from the perspective of the end-user [1].

The predominant method of creating E2E tests has relied heavily on human intervention, with developers manually assessing application features and using frameworks such as Selenium [2] to script the user scenarios. Efforts to automate E2E test generation have explored reinforcement learning (RL) [3] and model-based approaches [4], [5], [6]. More recently, the rise of large language models (LLMs) has spurred their application to various testing tasks. While LLMs have shown promise in generating unit tests for various applications [7], [8], [9], [10], [11] and mobile testing [12], [13], their application to web app E2E test generation remains an open area of research. Our recent work, FormNexus, has focused on automating web form testing [14], highlighting the potential for LLMs to enhance E2E testing methodologies.

In this paper, we formalize the notion of feature-driven E2E testing, including definitions of application features and a new metric called feature coverage for assessing E2E tests. Then,

we propose AUTOE2E, the first technique designed to generate feature-driven E2E tests autonomously. Our approach is centered on the ability to automatically infer features embedded within the web application and translate them into a sequence of user actions that form an E2E test scenario. We approximate potential features and employ a novel probabilistic scoring method that deduces the likelihood of a feature’s existence based on its observed frequency within the web application. In our probabilistic approximation, we leverage LLMs, such as GPT-4o [15] or CLAUDE 3 [16], which exhibit enhanced abilities to comprehend content within applications compared to traditional models.

A critical challenge we encountered was the absence of a suitable dataset for evaluating E2E test cases. To address this, we create a novel benchmark, called E2EBENCH comprising 8 open-source web applications. For each application, we extract all available features and employ instrumentation techniques to monitor front-end code, enabling us to track the actions performed on the web application during E2E tests. We establish a mapping between each application’s features and the corresponding sequences of actions performed. As part of the benchmark, we develop a tool capable of automatically monitoring E2E test suite execution and assessing its coverage across all existing features within each application.

In summary, this work makes the following contributions:

- A formal formulation of feature-driven E2E test case generation, providing a theoretical foundation for developing practical, automated E2E testing tools.
- A probabilistic method for automatically inferring features in a web application. Our feature inference system assesses the likelihood of feature existence based on the frequency of the observed user actions.
- A novel technique, called AUTOE2E, capable of autonomously generating feature-driven E2E test cases. Each test contains a sequence of actions to cover an inferred feature.
- E2EBENCH, a novel benchmark designed for automatic evaluation of E2E test suites, quantifying their effectiveness in feature coverage.

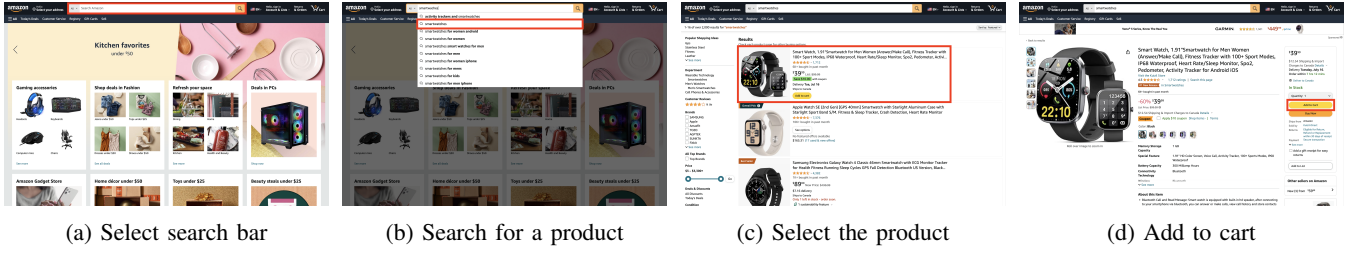


Fig. 1: Add to Cart feature on Amazon’s web application

Our evaluation results demonstrate that AUTOE2E achieves a feature coverage of **79%**, surpassing the best baseline, CRAWLJAX, by **558%**, and a significant **731%** improvement over the best LLM agent-based baseline, BROWSERGYM, highlighting the effectiveness of our methodology for automated E2E test generation. In addition to achieving superior coverage, AUTOE2E generates more complex test cases. Furthermore, the test cases ranked as more likely to exist by our likelihood estimation system exhibit a higher correspondence with actual features within the web applications, further validating the effectiveness of our approach.

## II. FEATURE-DRIVEN E2E TESTING

To illustrate the challenges and opportunities in E2E test generation, we utilize the Amazon web application [17] as a motivating example, with a few representative pages illustrated in Figure 1. Before deploying such apps, developers must conduct rigorous testing to guarantee proper functionality. Testing can occur at various levels, including unit tests, integration tests, and E2E tests. E2E tests, in particular, focus on executing specific features and functionalities that users will engage with from start to finish, ensuring that different features and scenarios within the app behave as expected. Within the context of E2E testing, we define:

**Definition 1** (User Operation). A user operation, denoted by  $U$ , is a sequence of user actions  $\{A_i\}$  (e.g., clicking, submitting forms). Each user operation can be classified into one of the following types:

- **Entity Operations:** These operations involve creating, reading, updating, or deleting data entities within the system. An entity operation is represented by a tuple  $(X, E, M, P)$ , where  $X$  signifies the CRUD action (Create, Read, Update, or Delete),  $E$  denotes the target entity or entity set,  $M$  is a boolean indicating whether the operation affects single or multiple entities, and  $P$  is a set of parameter values providing additional context or control.
- **Configuration Operations:** These operations modify system configurations or settings. They are represented by a tuple  $(C, P)$ , where  $C$  identifies the specific configuration being altered, and  $P$  specifies the new value.

On the Amazon web application, searching for a product is classified as an entity operation, specifically a read operation ( $X = \text{Read}$ ) on the Product entity set ( $E = \text{Product}$ ). As a search typically returns multiple results, the multiplicity is

true ( $M = \text{True}$ ), and the parameter set includes the search term entered by the user ( $P = \text{“search term”}$ ). It is important to distinguish between viewing a list of search results ( $M = \text{True}$ ) and viewing the details of a single product ( $M = \text{False}$ ), as these are distinct operations within the system.

Configuration operations are prevalent in various software systems. For example, logging into a system involves updating the Authentication configuration to reflect the user’s authorized status. Toggling a Dark Mode setting modifies the application’s visual appearance, while changing a Language configuration alters the language in which content is displayed to the user.

**Definition 2** (Application Feature). An application feature, denoted by  $\mathcal{F}$ , is characterized by the following properties:

- $\mathcal{F}$  can be realized through a finite ordered sequence of user operations (Definition 1), denoted as  $U_1 \rightarrow U_2 \rightarrow \dots \rightarrow U_n$ , where each  $U_i$  represents a distinct user operation.
- Each user operation  $U_i$  within the sequence is essential for achieving  $\mathcal{F}$ . Removing any  $U_i$  would result in an altered or unattainable outcome.
- The outcome  $O$  of executing the sequence  $U_1, U_2, \dots, U_n$  is visually presented to the user upon completion, signifying the successful realization of  $\mathcal{F}$ .
- The label  $L$  is an abstract, natural language description of the feature  $\mathcal{F}$ , independent of the parameters  $P_i$  in the operations  $U_i$ .

For instance, consider the common feature  $\mathcal{F}$  on the Amazon web application labeled as “Adding a Product to the Shopping Cart”. As illustrated in Figure 1, this feature entails a specific sequence of user operations  $U_i$ : (1) viewing (*Read*) a list of products, (2) viewing (*Read*) the details of a specific product, (3) Creating a cart item. This operation chain culminates in a visible change to the user’s cart, representing the outcome  $O$ .

Each operation in this sequence is crucial for the successful completion of  $\mathcal{F}$ . For instance, removing the search operation would prevent the user from finding the desired product, while omitting the product selection step would leave the system without a specific item to add to the cart. In contrast, operations not essential to achieving the feature, such as changing the application’s language, or background color should not be included in the defining sequence.

The abstract nature of  $\mathcal{F}$  labeling implies that it should not

be tied to specific parameters from the constituent operations. Whether the user searches for “shoes” or “electronics”, or chooses a particular brand or model, the fundamental feature of “Adding a Product to the Shopping Cart” remains unchanged.

**Manual Testing.** In practice, manual E2E test creation often relies on developers or quality assurance (QA) engineers exploring the application or utilizing design documents to extract features, functionalities, and user flows within the application, outlining key scenarios to test. After specifying the features, test scripts for those features are often written using specialized frameworks like Selenium [2], Cypress [18], or Playwright [19] to simulate user interactions within the app. These scripts typically involve executing specific actions on the app, such as clicking buttons, filling out forms, and navigating through pages. They also include assertions to verify that the application responds as expected at each stage.

**Test Generation.** Model-based testing has emerged as a prominent technique for automated E2E test generation in web applications [20], [21], [5], [6]. This approach involves systematically exploring the application’s state space by exercising available actions, such as clicking links, hovering over elements, and submitting forms. The resulting state transitions are then captured and integrated into a model representing the application’s behavior.

For instance, consider a model of the Amazon web application (illustrated in Figure 1). This model would encompass the states and action sequences necessary to add an item to the shopping cart. Furthermore, it would incorporate transitions facilitated by persistent navigation elements, such as the menu bar, enabling navigation between key pages like the landing page, login page, and shopping cart from any point within the application. Once constructed, this model serves as the foundation for automatic test case generation. Test cases are derived as sequences of actions extracted from the model’s transitions, with the objective of achieving comprehensive coverage based on predefined criteria, such as state coverage or transition coverage.

Listing 1 illustrates a model-based generated test case, showcasing a path within the Amazon web application.

```

1 test("generated test case 1", () => {
2   cy.visit("https://amazon.com");
3   cy.get("search-bar").type("smartwatch");
4   // assertion 1
5   cy.get("login").click();
6   // assertion 2
7   cy.get("home").click();
8   // assertion 3
9   cy.get("cart").click();
10  // assertion 4
11 });

```

Listing 1: Sample model-based generated test case

This Cypress-based test simulates user interactions such as searching for a product (“smartwatch”), navigating from the search results to the login page, returning to the landing page, and finally proceeding to the cart page. At each stage, assertions are made to verify the correctness of the application’s

state, ensuring the presence of expected elements, validating displayed content, and confirming successful navigation.

**Challenges in E2E Testing.** Both manual test creation and automated test generation present distinct challenges. Developer-written test cases, while potentially comprehensive, are often expensive and time-consuming to develop. Model-based techniques offer a partial solution by automating test case generation. However, the resulting tests may lack the coherence and relevance of human-written scenarios. This is evident in the sample model-based test case (Listing 1), which, despite covering several states and transitions, the sequence of actions performed is not relevant to a specific application feature (Definition 2). While developer-written tests typically follow meaningful features (e.g., adding an item to a cart), model-based tests prioritize coverage over coherence, leading to seemingly random sequences of actions. This trade-off between coverage and relevance highlights a key challenge in automated test generation. To more effectively evaluate the quality of generated E2E tests, we propose a new metric, feature coverage:

**Definition 3** (Feature Coverage). Let  $\mathcal{F} = \{f_1, f_2, \dots, f_m\}$  be the set of all features (Definition 2) in an application, and let  $\mathcal{T} = \{t_1, t_2, \dots, t_n\}$  be a set of test cases in a test suite designed to test these features. Define a relation  $R \subseteq \mathcal{T} \times \mathcal{F}$  where each  $(t_i, f_j) \in R$  indicates that test case  $t_i$  exercises feature  $f_j$ . Feature coverage  $C$  is then defined as the ratio of the number of unique features exercised by the test suite to the total number of features, given by:

$$C = \frac{|\{f_j \in \mathcal{F} \mid \exists t_i \in \mathcal{T} \text{ such that } (t_i, f_j) \in R\}|}{|\mathcal{F}|}$$

where each test case  $t_i$  targets exactly one feature  $f_j$  and it is permissible for multiple test cases to cover the same feature.

This metric shifts the focus from purely syntactical (code coverage) or structural coverage (such as state and transition coverage) to a more user-centric perspective, emphasizing the testing of distinct application features. Automating test case creation while achieving a high feature coverage necessitates a nuanced understanding of the application’s context and content, a capability that remains absent in existing automation techniques. To address this gap, a novel feature-driven E2E test generation approach is required. This approach must be capable of inferring the features existing within the app, connecting the available actions to those features, and generating the sequence of actions corresponding to the coverage of each feature as a test case.

### III. APPROACH

In this work, we introduce AUTOE2E, a novel approach for automatically generating semantically meaningful E2E tests, each targeting a distinct application feature (Definition 2), aiming to achieve high feature coverage (Definition 3). Our primary focus is addressing the challenge of inferring application features from the contextual information present in various application states. We propose a method that assesses

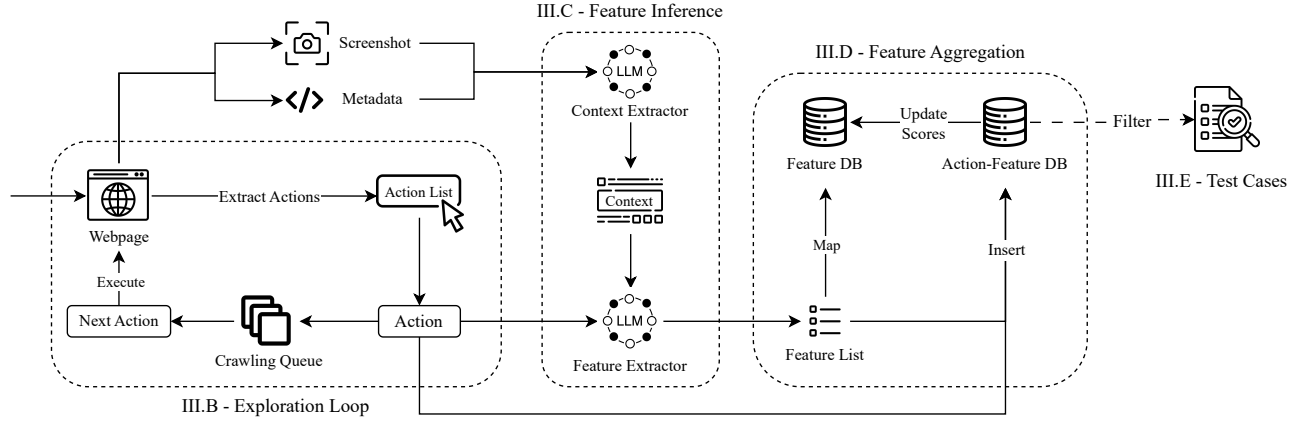


Fig. 2: Overview of our framework

the likelihood of features existing within an application based on observed user actions. By integrating this method with LLMs, we establish a workflow to identify the features present in the application and the corresponding sequences of actions required to trigger them. These action sequences are then transformed into comprehensive test cases for the application. The architecture for AUTOE2E is illustrated in Figure 2.

#### A. Feature Inference Modeling

We first formalize the feature inference task in order to leverage LLMs more effectively. A web app typically consists of various application states and transitions between them:

**Definition 4** (Application States and Transitions). An application state  $S_i$  represents a snapshot of the web app at a particular moment, characterized by the runtime values of relevant variables, as well as the dynamic structure and content as rendered in the browser. A transition  $A_i$  initiated by a user action (e.g., clicking a button, submitting a form) can cause a state change, e.g., from  $S_1$  to  $S_2$ .

Figure 1 provides a visual representation of this concept, where each image depicts a distinct state within the Amazon web application. Furthermore, the actions available within these images represent potential transitions to other states.

Given a web app with  $K$  states,  $\{S_1, S_2, \dots, S_K\}$ , and  $M$  potential features,  $\{\mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_M\}$ , the task of inferring features becomes that of constructing a generative model to estimate the following distribution:

$$p(\mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_M | S_1, S_2, \dots, S_K) = p(\mathbf{F} | \mathbf{S}) \quad (1)$$

In this formulation, the probability distribution  $p(\mathbf{F} | \mathbf{S})$  represents the likelihood of a feature set  $\mathbf{F}$  being present within an application, given the information observed in the set of states  $\mathbf{S}$ . A generative model could learn this distribution from data and subsequently generate a feature set  $\mathbf{F}$  that maximizes this probability.

This approach aligns with how humans typically extract features and design test scenarios. Users explore the application’s interface, inferring available features based on observations. For instance, when presented with the application states depicted in Figure 1, a human can intuitively infer the presence of features such as searching for products, purchasing products, viewing account details, and reviewing order details. These inferences are based on the visual cues and interactive elements present in the observed states, leading to a higher likelihood of these features being available within the Amazon web application.

However, constructing a generative model capable of accurately estimating  $p(\mathbf{F} | \mathbf{S})$  is a challenging task due to the vast complexity of both the feature and state space. To address this challenge, we introduce a simplified model that enhances tractability while retaining the essence of the generative approach.

1) *Feature Independence*: Given complete access to all states  $\mathbf{S}$  within an application, the inference of individual features  $\mathcal{F}_i$  can be considered independent. While certain features may often imply the existence of others (e.g., an “add to cart” feature suggests a “remove from cart” functionality), knowledge of the complete state space allows for direct observation and inference. For instance, the presence of a “remove” button on the Cart page confirms the existence of the “remove from cart” feature, independent of knowledge about the “add to cart” feature. Leveraging this feature independence given  $\mathbf{S}$ , the joint probability distribution in Equation 1 simplifies to:

$$p(\mathbf{F} | \mathbf{S}) = p(\mathcal{F}_1 | \mathbf{S}) p(\mathcal{F}_2 | \mathbf{S}) \dots p(\mathcal{F}_M | \mathbf{S})$$

Based on this notion, instead of attempting to infer all features simultaneously, we can leverage the conditional independence and generate features individually using the distribution  $p(\mathcal{F}_i | \mathbf{S})$ . By sampling the top  $M$  generated values from this distribution, we can effectively identify the most probable existing features within the application.



2) *Action-Centric Feature Inference*: Feature determination in web applications can typically be accomplished through an action-centric lens, focusing on the available actions on a page rather than the specific content. As illustrated in Figure 1d, the “Add to Cart” feature on a product page is discernible solely from the context of the page and the presence of the corresponding action, regardless of the product’s details. Having a general context for the page is particularly valuable when dealing with actions that have ambiguous descriptions, such as a “continue” button. In such cases, the context of the current page (e.g., checkout page) aids in clarifying the intended purpose of the action. By adopting this action-centric perspective, we shift the focus from analyzing the entirety of the application’s state information to a more targeted examination of the actions and their associated contextual information. Formally, we can express this as:

$$p(\mathcal{F}|\mathbf{S}) = p(\mathcal{F}|A_{1,1}, A_{1,2}, \dots, A_{K,n_K}) \quad (2)$$

where  $A_{i,j}$  and  $n_i$  represent the  $j$ -th action and the number of actions on state  $S_i$  respectively.

3) *Sequential Action Chains*: The features within web apps are designed to be executed through a sequential chain of actions, as can be observed in Definitions 1 and 2. For a given feature  $\mathcal{F}$ , there exists an ordered sequence of actions  $A_1, A_2, \dots, A_N$  that leads to its execution:

$$\mathcal{F} : A_1 \rightarrow A_2 \rightarrow \dots \rightarrow A_N$$

Crucially, a feature should be derivable solely from its action chain. For example, when adding an item to a shopping cart, the presence of actions for rating or commenting on products is irrelevant. This allows us to eliminate unrelated actions from the context of Equation 2, focusing solely on the actions  $\{A_i\}$  relevant to the feature:

$$p(\mathcal{F}|\mathbf{S}) = p(\mathcal{F}|A_1, A_2, \dots, A_N)$$

Then we can use Bayes’ theorem and chain rule of probabilities to simplify this to:

$$\begin{aligned} &= \frac{p(\mathcal{F}, A_1, A_2, \dots, A_N)}{p(A_1, A_2, \dots, A_N)} \\ &= \frac{p(\mathcal{F})p(A_1|\mathcal{F})p(A_2|A_1, \mathcal{F}) \dots p(A_N|A_1, \dots, \mathcal{F})}{p(A_1, A_2, \dots, A_N)} \end{aligned} \quad (3)$$

4) *Dependence of Subsequent Actions*: The sequential nature of user interactions within web applications can be modeled by recognizing the dual dependency of each action: an action  $A_i$  depends on the immediately preceding action  $A_{i-1}$  and the specific feature  $\mathcal{F}$  the user intends to execute. This dependency allows us to estimate the most probable subsequent actions based on the current action and feature.

Consider the scenario illustrated in Figure 1c, where a user clicks on a product within a list of search results. If the user intends to purchase the product (feature  $\mathcal{F}$ ), we can predict that the next likely action would be clicking on “Add to Cart.” Conversely, if the user intends to rate the product (a different feature  $\mathcal{F}$ ), the most probable next action would be clicking

on a rating value. Crucially, this estimation of the next action relies solely on the current action and the intended feature, regardless of the user’s prior interactions and the specific path taken to reach the current state. Whether the user arrived at the product page through searching, browsing categories, or any other means is irrelevant. Consequently, based on this property, we can simplify Equation 3:

$$\begin{aligned} p(\mathcal{F}|\mathbf{S}) &= \alpha(\mathbf{A})p(\mathcal{F})p(A_1|\mathcal{F}) \prod_{i=2}^N p(A_i|A_{i-1}, \mathcal{F}) \\ p(\mathcal{F})p(A_1|\mathcal{F}) &= p(A_1)p(\mathcal{F}|A_1) \\ p(A_i|A_{i-1}, \mathcal{F}) &= \frac{p(\mathcal{F}|A_i, A_{i-1})p(A_i|A_{i-1})}{p(\mathcal{F}|A_{i-1})} \\ \rightarrow p(\mathcal{F}|\mathbf{S}) &= \beta(\mathbf{A})p(\mathcal{F}|A_1) \prod_{i=2}^N \frac{p(\mathcal{F}|A_i, A_{i-1})}{p(\mathcal{F}|A_{i-1})} \end{aligned} \quad (4)$$

In the presented equations,  $\alpha(\mathbf{A})$  and  $\beta(\mathbf{A})$  denote generalized action probabilities. These terms are not specific to any particular application or feature but rather capture the inherent likelihood of observing certain actions across all the different web interactions. For instance, these terms might encapsulate the probability of encountering a “Login” action in any web application, essentially providing a baseline expectation for the occurrence of actions. Since our goal is to maximize the distribution over  $\mathcal{F}$ , these functions become irrelevant. Therefore, we can further simplify the task of feature inference into the following equation:

$$\begin{aligned} \mathcal{F} &= \arg \max_{\mathcal{F}} \sum_{i=1}^N \left( \log(p(\mathcal{F} | A_i, A_{i-1})) \right. \\ &\quad \left. - \log(p(\mathcal{F} | A_{i-1})) \right) \end{aligned} \quad (5)$$

This result has an intuitive interpretation. If we are predicting the existence of a certain feature based on an action, we should observe further evidence supporting that feature after performing the action.

Equation 5 provides a flexible foundation for various implementations, as it applies broadly to web applications and is not constrained by any specific implementation details in its derivation. In this work, we implement our method based on this equation by utilizing LLMs to estimate the distributions  $p(\mathcal{F}|A_i)$  and  $p(\mathcal{F}|A_i, A_{i-1})$ . This is achieved by feeding  $\{A_i\}$  and  $\{A_{i-1}, A_i\}$  as context to the LLM, respectively, and prompting it to infer features based on this context. We then aggregate the generated results to assess the likelihood of existing features within the application. The remaining sections of our approach will detail how we leverage LLMs and aggregate their outputs to infer both features and their corresponding chains of actions.

## B. Exploration Loop

As illustrated in Figure 2, AUTOE2E operates on an exploration loop paradigm, interfacing with the target web application, capturing user actions, and systematically queuing them

for execution. Each executed action potentially reveals new application states, driving exploration until no further novel states are discovered or a timeout occurs. AUTOE2E employs a breadth-first search (BFS) strategy for state exploration, queuing, and recursively crawling neighboring states from the current state. Feature extraction inferences are performed concurrently during state visits, feeding the extracted actions into the rest of the workflow to infer features.

### C. Feature Inference

As the exploration loop discovers new states and extracts their associated actions, AUTOE2E concurrently analyzes each action to infer the specific features with which it interacts. For the following sections, imagine our exploration has led us to state  $S_i$  via the action sequence  $A_1 \rightarrow \dots \rightarrow A_{i-1}$ . In this state, we observe a set of available actions  $\mathbf{A}_i = \{A_{i1}, A_{i2}, \dots, A_{in}\}$ , where  $A_{ij}$  denotes the  $j$ th action on  $S_i$ .

1) *State Context Extraction*: As discussed in Equation 2, actions are represented alongside their corresponding context—the high-level purpose of the page where the action occurs. This contextual information is essential for accurate feature inference, particularly when the content associated with an action is ambiguous (e.g., a button labeled “Continue”) and requires further disambiguation.

Identifying page context requires multiple data sources. The application’s description and category provide high-level information, while the page’s content (text, HTML, images) offers more specific details. We prioritize image-based analysis over HTML due to HTML’s verbosity. Additionally, the history of actions leading to the current page provides further contextual clues. Consequently, our context extraction employs a multi-modal approach, incorporating a screenshot of the  $S_i$  rendered in a browser, a description of the entire application, and the most immediate action  $A_{i-1}$  leading to  $S_i$ . As an example, the following is the extracted context for the page in Figure 1c: *A webpage displaying search results for a product query, allowing users to browse and filter options for purchasing.*

2) *Feature Extraction*: Following the extraction of contextual information from  $S_i$ , we can now utilize that context in conjunction with the actions present on the state,  $\mathbf{A}_i$ , and prompt an LLM to predict possible features connected to each action. We follow the result in Equation 5 derived in Section III-A to infer the features by querying the LLM twice. The first prompt requests the LLM to generate features based on the individual actions in  $S_i$ , corresponding to  $p(\mathcal{F}|A_{ij} \in \mathbf{A}_i)$ . The second prompt asks the LLM to generate features based on the actions in  $S_i$  plus the most recent action that led to the current state, representing  $p(\mathcal{F}|A_{ij} \in \mathbf{A}_i, A_{i-1})$ . To illustrate this process, consider Figure 1d. For the first prompt, the LLM would be asked to generate features solely based on the “Add to Cart” button. However, for the second prompt, the LLM would consider both the “Add to Cart” button and the preceding action that led to this page, which is clicking the product link in Figure 1c.

Estimating the probability of a result generated by an LLM, particularly in proprietary models, is not always feasible. This

limitation hinders the direct calculation of  $p(\mathcal{F}|A_{ij} \in \mathbf{A}_i)$  and  $p(\mathcal{F}|A_{ij} \in \mathbf{A}_i, A_{i-1})$  in Equation 5. To address this, we incorporate Chain of Thought (CoT) prompting in the LLM prompt, asking it to generate a list of  $R$  features ordered by their perceived probability of existence. CoT prompting encourages the LLM to provide more reasoned and reliable responses, increasing the validity of the feature ordering. Having access to this ordering, we employ a geometric distribution to estimate the probability of an inferred feature based on its rank:

$$\begin{aligned} \text{rank\_score}(r|r \leq R) &= \log(p(\mathcal{F} \text{ is ranked } r)) \\ &= \log((1-p)^{r-1}p) = (r-1)\log(1-p) + \log(p) \end{aligned} \quad (6)$$

where  $p$  is a manually set parameter. A  $p$  value close to 1 creates a significant difference in probability between the top-ranked item and the rest, while a  $p$  value near 0 assigns nearly equal probability to each rank. As there is often more than one feature associated with an action, we seek a balance that allows our model to recognize these multiple features while still distinguishing between higher and lower-ranked items. Therefore, we set  $p$  to 0.5.

Since the list of our inferred features is limited to the top  $R$ , we must also take into account the probability that a feature is not in the top  $R$ , but appears at a certain rank if we extend our list. For every feature that does not appear in the top  $R$  inferred features, we use a constant score:

$$\text{rank\_score}(r|r > R) = R\log(1-p) + \log(p) \quad (7)$$

This constant score is used in the aggregation phase.

### D. Feature Aggregation

During the exploration process, upon encountering a new state, we extract potential features associated with each action within that state. Importantly, these feature inferences are conducted in isolation. Features derived using the current action context  $\{A_{ij}\}$  are independent of those inferred using both the current and preceding actions  $\{A_{ij}, A_{i-1}\}$ . Furthermore, both sets of newly inferred features are initially disconnected from any previously generated features. To reconcile these disparate inferences, we enter the aggregation phase. To manage the aggregation, we employ two interconnected databases.

1) *Feature Database (FD)*: The Feature Database (FD) is a vector database storing a global list of discovered features. Each entry in FD contains a label, i.e., a textual description of the feature, its corresponding embedding, and a confidence score derived from Equation 5. This score reflects the likelihood of the feature’s existence within the application.

2) *Action-Feature Database (AFD)*: FD is complemented by the Action-Feature Database (AFD). Each row in AFD associates an action within a state with its corresponding inferred features from FD. Additionally, it records a rank score calculated using Equation 6 and indicates whether the inference utilized solely the current action  $A_{ij}$  or both the current state’s action and preceding action  $A_{ij}, A_{i-1}$  as context.

The information stored in FD and AFD is dynamic, evolving as the application is explored. These databases continuously interact, influencing each other to update the overall confidence scores of features. This iterative refinement ensures that the feature model remains accurate and comprehensive as new observations are made.

3) *Mapping inferred features to FD*: To update feature scores, we first map the newly inferred features to existing entries in the FD. This is achieved by leveraging the textual embeddings of feature descriptions stored in FD. We query FD using the embedding of an inferred feature and retrieve the top results based on cosine similarity. This metric identifies features in FD that are semantically closest to the inferred feature. By establishing these connections, we can then update the confidence scores of the corresponding features in FD, incorporating the information gained from the new observations.

While querying the FD using textual embeddings can identify semantically similar features, it does not guarantee a perfect match. To ensure accuracy, we introduce an additional validation step using an LLM. This step involves querying the LLM to assess whether the inferred feature aligns with any of the similar feature descriptions retrieved from FD. The LLM acts as a semantic arbiter, determining if any of the retrieved descriptions truly describe the same feature. Listing 2 provides a concrete example of this process, showcasing how the LLM helps refine the initial matches obtained from FD for an “Add to Cart” action’s feature in Figure 1d.

```

1 // A feature inferred for an action
2 "Add a product to the cart"
3
4 // FD retrieval results
5 "Adding to the cart"           // exact match
6 "View product details"       // mismatch
7 "Remove item from the cart"   // mismatch
8 "Explore product categories"  // mismatch
9 "Navigate to homepage"       // mismatch

```

Listing 2: Example for FD querying

If no match is found in FD, this indicates that the feature is a novel observation within the application. In such cases, the feature data, along with an initial confidence score of 0, is inserted into FD. Following the matching (or insertion) of a feature in FD, we then create a corresponding entry in the AFD. This entry includes a pointer to the relevant FD entry, establishing the connection between the two databases. With these updates in place, we can then proceed to refine the confidence score of the feature in FD, incorporating the information gained from the new observation.

4) *Updating the Scores*: For each action  $A_{ij}$ , we have inferred a corresponding set of features  $\{\mathcal{F}_{j1}, \mathcal{F}_{j2}, \dots\}$ , as detailed in Section III-C. To update the feature scores based on these new observations, we employ the following formula, derived from Equations 5 and 6:

$$\text{score\_update}(F_{jk}) = \text{rank\_score}(F_{jk}|A_{ij}, A_{i-1}) - \text{rank\_score}(F_{jk}|A_{i-1}) \quad (8)$$

This formula calculates the change in the score for feature  $F_{jk}$  (the  $k$ th feature associated with action  $A_{ij}$ )

based on the observed action sequence. The first term,  $\text{rank\_score}(F_{jk}|A_{ij}, A_{i-1})$ , represents the probability of the feature given both the current and preceding actions. The second term,  $\text{rank\_score}(F_{jk}|A_{i-1})$ , represents the probability of the feature given only the preceding action. The difference between these two terms reveals the incremental impact of action  $A_{ij}$  on our confidence in the existence of feature  $F_{jk}$ .

To address cases where feature  $F_{jk}$  exists in the LLM’s inference for  $(A_{ij}, A_{i-1})$  but not for  $A_{i-1}$ , we substitute the  $\text{rank\_score}(F_{jk}|A_{i-1})$  term in Equation 8 with the constant ranking score from Equation 7. Using the *score\_update* formula, we then locate  $F_{jk}$ ’s corresponding entry in FD and update its score. This process is repeated throughout the app exploration.

#### E. Generating Test Cases

Following the exploration phase, we sort the features within the FD. From this sorted list, we filter and retain the top-scoring features as those identified within the application. This filtering process involves selecting a lower bound cutoff, determined by Equation 6 as  $\log(p)$ . This choice of cutoff is motivated by the possibility of single-action features, which could be correctly predicted by the LLM with a rank of 1 and a corresponding score of  $\log(p)$ . By setting the cutoff at this value, we ensure the inclusion of such features in our generated tests.

After filtering, we then extract the corresponding chains of actions from the AFD and translate them into test cases. This process leverages the accumulated evidence gathered during exploration to identify the most likely features and their associated action sequences, ultimately generating test cases that semantically cover the application features.

#### F. Implementation

AUTOE2E is implemented in Python using the LangChain framework [22], offering flexibility in LLM selection. For our evaluations, we opted for CLAUDE 3.5 SONNET due to its recognized performance among the most advanced LLMs. Textual embeddings for feature descriptions and FD utilize the ADA architecture [23], and the E2E tests are generated in Selenium [2]. Our databases, FD and AFD, are hosted on MongoDB Atlas<sup>1</sup>, which provides a vector search functionality well-suited for querying the feature description embeddings.

### IV. BENCHMARK CONSTRUCTION

Assessing E2E test cases presents a significant challenge in software testing automation. To the best of our knowledge, there is an absence of datasets for E2E test case evaluation that hinders further progress in this area of study. To address this gap, we embarked on the creation of such a benchmark, called E2EBENCH. This section elucidates the complexities inherent in dataset creation and outlines the methodological steps undertaken to achieve this goal.

**Benchmark Construction Challenges.** The lack of a benchmark for assessing Feature Coverage in E2E test cases arises

<sup>1</sup><https://www.mongodb.com/>

from the inherent challenges of evaluating this metric. As evident in Definition 3, this evaluation hinges on two critical steps: identifying the features within an application and mapping E2E test cases to these features to measure coverage.

The first challenge in this process is the identification and quantification of features within a web application. Features are often difficult to define and vary widely depending on the context, making their extraction difficult and subjective. Once features are identified, the second challenge is determining which specific feature a test case targets. This becomes particularly complex in large-scale applications like Amazon (illustrated in Figure 1), where a feature such as “Viewing the product’s details” can be accessed through multiple pathways—whether by clicking on a product from the landing page, filtering through categories, or using the search function. As the number of available actions increases, so does the complexity of tracking these diverse paths. For example, the Amazon web application with millions of products could have numerous valid E2E test cases for viewing product details, each following a distinct path. An effective evaluation tool must be capable of accounting for and managing this inherent path redundancy.

**Feature Identification.** Regarding the first challenge, we have established a precise and formal definition of a feature (Definition 2) in this work, providing a foundation for systematic and objective identification of features within software applications.

**Instrumentation for Feature Mapping.** To address the second challenge in mapping test cases to features, we need to track user (i.e., E2E test) interactions within web applications. To automate this tracking, we select open-source web applications for our benchmark, enabling us to directly instrument their code. The instrumentation adds logging mechanisms to capture every user action within the app, encompassing clicks, hovers, and various input types (selects, texts, radio buttons, checkboxes). Each unique action component within the app generates a distinct log message containing a unique identifier. To handle path redundancy, we leverage the inherent modularity of web applications. For example, consider an open-source e-commerce application: products are typically rendered using a collection of `Product` components. By leveraging this component-based structure, we can instrument component code to track the interactions by the E2E test cases. This approach effectively circumvents the challenge of near-infinite pathways, as each action within a path is implemented within a discrete component, regardless of its frequency of occurrence.

**Feature Grammar Extraction.** With the tracking system established, we then map the chain of logs to specific features within each subject application. This process involves identifying the available features for each subject in the benchmark and executing sequences of actions to trigger them while tracking the corresponding logs. These logs are subsequently transformed into a *feature grammar*, which serves as a ground-truth reference for assessing E2E test cases. To illustrate, a

sample feature grammar for purchasing a product in an e-commerce application would be the following:

```
"c1-product-element" "c10-add-to-cart"
"c22-cart-icon" "c12-checkout-button"
("t5-credit-number" | "t16-credit-date" |
 "t23-credit-cvv")+ "c35-complete-purchase"
```

Listing 3: Feature Grammar for “product purchase”

In this feature grammar, the sequence of actions begins by clicking on the product, followed by adding it to the cart and navigating to the cart page. Subsequently, the checkout button is clicked, leading to the purchase page where credit card information is required. The input fields for credit card details can be filled in any order, represented by the logical OR operator. Additionally, the information in the text boxes can be modified multiple times, denoted by the “+” sign indicating repetition. Finally, the purchase is completed by clicking a confirmation button.

The extraction of feature grammar for each subject was conducted independently by each author, followed by a collaborative discussion to consolidate the findings. This iterative process ensured a comprehensive identification of features, with consensus reached on both the features themselves and any alternative paths to achieve the same functionality. The resulting grammars then serve as a basis for evaluating the coverage of a test suite over all of the functionalities.

**Benchmark Subjects.** A list of the apps in our benchmark is available in Table I. These applications, most of which have been previously used in web testing research [24], [25], [26], [27], encompass a diverse range of categories, including bug tracking (*MantisBT*), e-commerce (*Saleor*), and translation management (*EverTraduora*).

**Automatic Coverage Evaluation.** During test case execution, our benchmark continuously monitors the performed actions, trying to map the sequence to one of the identified functionalities within the application. This allows us to assess the coverage of an E2E test suite across all the distinct features in our benchmark subjects. The calculation of feature coverage is performed completely automatically based on the grammar extracted in the previous phase, resulting in an objective measurement of the Feature Coverage (Definition 3).

## V. EVALUATION

We have framed the following research questions to measure the effectiveness of AUTOE2E:

- **RQ1:** How effective is AUTOE2E in generating feature-driven E2E tests?
- **RQ2:** How accurate is the feature inference of AUTOE2E?
- **RQ3:** How does AUTOE2E compare to other state-of-the-art techniques?

**Process.** We use E2EBENCH to evaluate the efficacy of AUTOE2E and other methods in achieving feature coverage. For running our experiments, we set the temperature parameter of the LLMs to 0 to produce the same response every time.

TABLE I: Benchmark Subjects

App Name	Category	Features	LOC
PetClinic	Health	23	51K
Conduit	Blog	17	53K
Taskcafe	Task Manager	32	67K
Dimeshift	Expense Tracker	21	10K
MantisBT	Bug Tracker	27	118K
EverTraduora	Translation Manager	41	25K
Saleor Storefront	E-Commerce	13	58K
Saleor Dashboard	E-Commerce Admin	130	1.1M

**Baselines.** To the best of our knowledge, no existing technique directly addresses feature-driven E2E test generation. Recent advances in LLM-based agents such as **WEBCANVAS** [28] and **BROWSERGYM** [29] have demonstrated the potential to navigate web applications and execute user instructions. However, adapting these agents for E2E test generation poses several challenges. First, they often require pre-existing feature extraction, which is the primary focus of our approach. Second, these agents may struggle to interpret abstract task descriptions and instead require concrete, detailed instructions for execution.

To assess the effectiveness of our proposed methodology, we evaluate AUTOE2E against these specific-purpose agents, **WEBCANVAS** and **BROWSERGYM**, both of which utilize GPT-4o as their underlying LLM. We also compare AUTOE2E with more generalized agents: **AutoGPT** [30], which employs GPT-4o for both task planning and execution, and **OpenDevin** [31], a code-focused agent that utilizes the **CLAUDE 3** LLM. Additionally, we benchmark against a model-based technique, represented by **CRAWLJAX**’s test generation module [20], [5]. All agents are instructed to navigate web applications and generate end-to-end (E2E) test cases.

#### A. Effectiveness (RQ1)

The primary objective of our method is to maximize the extent of *Feature Coverage* as defined in Definition 3. As described in Section III, AUTOE2E generates a set of E2E test cases. We evaluate the generated tests using the methodology detailed in section IV.

AUTOE2E demonstrates the ability to generate test cases that cover an average of **79%** of features across the applications in the E2EBENCH benchmark. Considering all features across all applications, AUTOE2E achieves a total Feature Coverage of 72%. A more granular breakdown of the feature coverage for each app is presented in the *Recall* column of Table II. These results underscore the effectiveness of AUTOE2E in effectively inferring features and generating corresponding test cases across a diverse range of applications.

#### B. Feature Inference (RQ2)

As described in Section III, AUTOE2E infers a list of features along with their scores by the end of the inference phase. Table II provides statistics for both correct and incorrect predictions within the inferred feature list. AUTOE2E achieves a total precision of 0.55, indicating that 55% of the generated

test cases were correct based on the evaluation metrics. The total recall of the model, which is the same as the average Feature Coverage across all features in all applications, reaches 0.72. This means that our generated test cases successfully cover 72% of the total features present. The precision and recall values result in an F1 score of 0.62, demonstrating a reasonable balance between AUTOE2E’s ability to identify relevant test cases and its ability to capture the full spectrum of application features.

TABLE II: Inference Statistics

App Name	Total	Correct	Precision	Recall	F1
PetClinic	29	19	0.66	0.83	0.73
Conduit	19	14	0.67	0.82	0.76
Taskcafe	50	22	0.42	0.66	0.51
Dimeshift	40	17	0.42	0.81	0.56
MantisBT	45	22	0.47	0.75	0.58
EverTraduora	69	30	0.42	0.73	0.54
Storefront	21	13	0.62	1.0	0.76
Dashboard	123	85	0.69	0.65	0.67
Total	396	222	0.55	0.72	0.62

However, the effectiveness of our approach is not solely determined by the number of correct features. As detailed in Section III, we employ a scoring and filtering system to prioritize feature generation and test case creation. If this system functions as intended, features with higher scores should be more likely to correspond to actual features within the application.

Figure 3 illustrates the coverage of top- $k$  features in relation to the ratio of  $k$  to the total number of features present in the application, referred to as the *Rank Ratio*. The left half of this figure demonstrates the coverage against the rank ratio for AUTOE2E, while the lower half displays the moving average of coverage for AUTOE2E and the baseline methods. In an ideal scenario, where all top- $k$  inferred features are accurately identified, the plotted line would follow a 45-degree trajectory. The figure reveals the extent of coverage achieved as the feature list expands to encompass more candidates.

As evident in the figure, AUTOE2E generally follows the expected 45-degree line up to a rank ratio of approximately 0.75. This indicates that, on average, for an application with  $N$  features, the first  $0.75N$  generated test cases are almost all correct. However, a divergence is observed beyond this point, suggesting that while the remaining generated test cases may be correct in certain instances, they do not consistently cover features with the same level of accuracy.

#### C. Comparison (RQ3)

We evaluated the generated tests by each baseline against E2EBENCH as described in section IV. The Feature Coverage results for different subjects and baselines are presented in Figure 4. AUTOE2E achieves an average Feature Coverage rate of **79%**, significantly outperforming **CRAWLJAX** (12.0%), **WEBCANVAS** (0%), **BROWSERGYM** (9.5%), **AUTOGPT** (6.1%), and **OPENDEVIN** (7.9%). Notably, AUTOE2E surpasses the

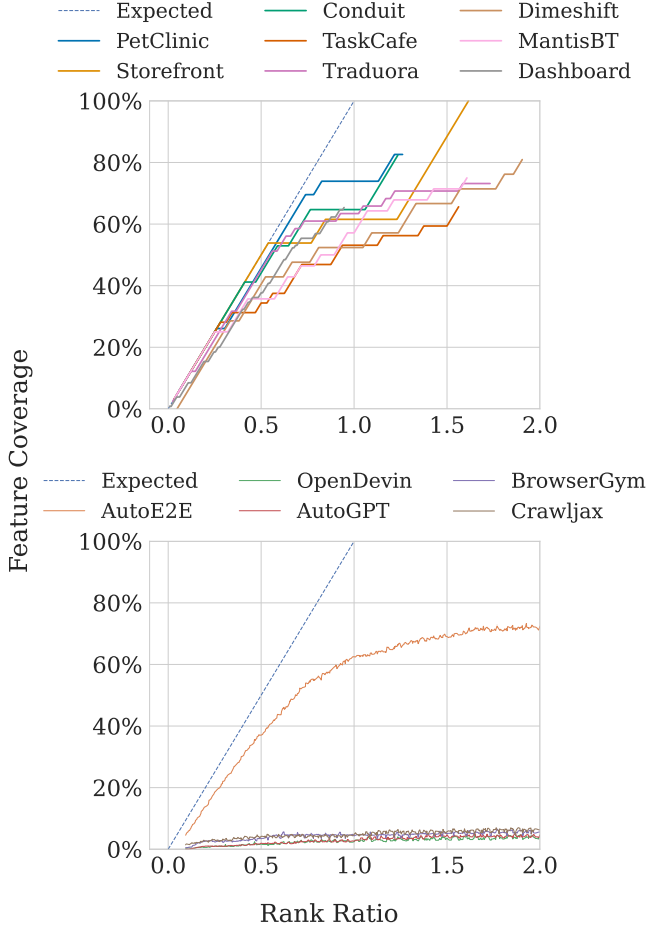


Fig. 3: Feature Coverage vs. rank ratio

next best performing tool, CRAWLJAX, by **558%**, and best agent-based tool, BROWSERGYM, by **731%**. WEBCANVAS has not been included in Figure 4 since it did not generate any test cases.

The F1 scores for the generated test cases were 0.62 for CRAWLJAX, 0.11 for BROWSERGYM, 0.08 for AUTOGPT, and 0.08 for OPENDEVIN, compared to 0.62 for AUTOE2E.

1) *Feature Complexity*: While the coverage rate treats all features equally, it is important to acknowledge that features vary in complexity. Some features within an application necessitate longer chains of ordered actions to be successfully executed. Different tools may encounter difficulties with extended action chains due to the challenge of maintaining context over longer periods. Table III provides a statistical analysis of action chain complexity across different tools.

The average length of feature chains in the benchmark is 3.4 actions. Notably, AUTOE2E demonstrates proficiency in handling longer chains, averaging 3.8 actions per feature, compared to 2.9 for CRAWLJAX, 1.4 for BROWSERGYM, 1.2 for AUTOGPT, and 1.7 for OPENDEVIN.

2) *Test Case Count*: The number of test cases generated varies across different tools, as each tool employs its own criteria for determining the appropriate quantity. This raises

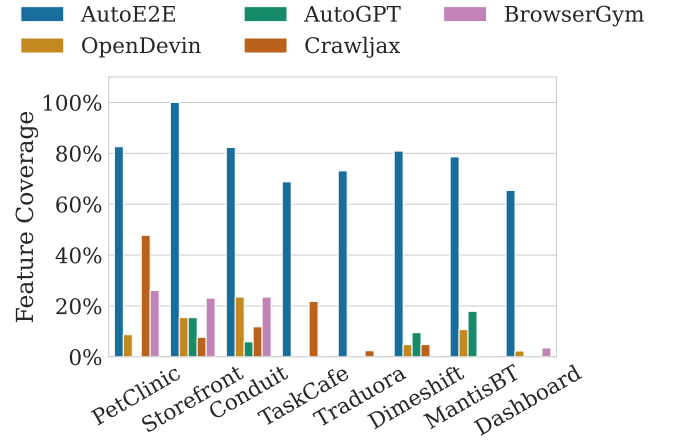


Fig. 4: Feature Coverage of different methods on the subjects

TABLE III: Feature action chain length

Tool Name	Min	Max	Average	Median
AUTOE2E	1	7	3.8	4
CRAWLJAX	1	7	2.9	3
WEBCANVAS	0	0	0.0	0
BROWSERGYM	1	3	2.0	2
AUTOGPT	1	3	1.2	1
OPENDEVIN	1	4	1.7	1
All Features	1	7	3.4	4

the question of whether forcing a tool to generate more test cases would necessarily lead to increased feature coverage. The lower half of Figure 3 sheds light on this by illustrating the moving average of Feature Coverage for the methods across all subjects. Notably, the baselines reach a plateau in coverage at a relatively low rate, suggesting that simply increasing the number of test cases does not guarantee a corresponding increase in feature coverage.

## VI. DISCUSSION

The framework developed in this work, particularly the findings from Section III-A, is platform- and implementation-agnostic. This means the underlying principles can be extended to generate test cases for other platforms, such as mobile applications. Furthermore, this implementation independence allows for significant potential improvements in performance and cost-effectiveness in future iterations, as the framework is not rigidly tied to any specific technology or tool. Additionally, the introduction of E2EBENCH provides a standardized and automated means of evaluating E2E test case generation techniques. This benchmark fills a crucial gap in the research community, enabling more rigorous comparison and development of novel approaches in the field of automated E2E test generation.

**Limitations.** Despite its strengths, our approach has limitations. Currently, AUTOE2E focuses on a one-to-one mapping between test cases and features, whereas real-world scenarios



often require multiple test cases per feature to assess diverse interactions. Additionally, our implementation generates assertions using the entire state after each action, which may not always be the most meaningful or targeted approach.

**Threats to Validity.** It is important to acknowledge potential threats to the validity of our findings and the steps taken to mitigate them. One such threat lies in the representativeness of the web applications selected for E2EBENCH. To address this, we have carefully curated a diverse set of applications spanning a wide range of categories. This diversity aims to enhance the generalizability of our results and reduce potential biases associated with specific application types.

Another potential threat stems from the inherent subjectivity in defining and extracting features. To mitigate this, we have established a precise and formal definition of a feature (Definition 2). Furthermore, we employed multiple authors in the feature identification process during benchmark construction, promoting a more objective and comprehensive perspective.

Finally, the use of different LLMs in the evaluation process could pose a threat to validity. While our goal was to maintain consistency by utilizing the same LLMs across all methods, certain constraints necessitated variations. Specifically, among the selected agents and baselines, only OPENDEVIN supported CLAUDE 3.5 SONNET, while others (AUTOGPT, WEBCANVAS, BROWSERGYM) lacked support for Anthropic models at the time of evaluation, requiring the use of GPT-4O. However, it is important to note that both CLAUDE 3.5 SONNET and GPT-4O exhibit very similar performance in reported benchmarks [32], [33]. Therefore, we anticipate that these variations in LLM usage would not significantly impact the overall performance trends observed in our evaluation.

## VII. RELATED WORK

**Web Navigation Agents.** In the academic context, extensive research has been conducted on automating the execution of tasks defined in natural language [34], [35], [36], [37], [38], [39], [40]. The aim is for an agent to determine and execute a sequence of actions within a web application that fulfills these instructions. This research is divided into two principal categories: traditional methods utilizing Reinforcement Learning (RL) agents and newer approaches that focus on LLMs.

The traditional techniques involve deep learning models to translate natural language instructions into embeddings [34], [35], [36], [37], [39], [40]. These embeddings are then mapped to specific actions on the web page. Variations in these methods are seen in the architecture of the deep learning models, the strategies of the RL policies, and how actions are modeled. Additionally, some methods might include or omit certain components. Techniques in this category often incorporate demonstrations by expert users [39], [40], use heuristics based on human-designed systems [41], or confine the agent’s actions to a predetermined set [37], [38].

Contrastingly, newer methods leveraging LLMs typically forgo the RL training phase. These models delegate the decision-making about subsequent actions to the LLMs [42],

[43], [44], [45]. These newer agents have been proposed by academic research, open-source communities [22], [30], and large-scale companies [15].

**LLM-based Testing.** Recent studies have increasingly focused on using LLMs to automate software and web testing processes. Some research concentrates on software unit test generation [46], [47], [48]. Some studies focus on accessibility testing, utilizing LLMs to identify and address accessibility issues [49]. Other studies have directed their attention towards GUI testing, including software GUI testing [50] and mobile application GUI testing [12], [13], [51]. Our recent work has focused on automated web form testing, where LLMs simulate user interactions and validate form functionalities through constraint-based testing [14].

**E2E Test Generation.** A specific but less extensive area of research focuses on generating end-to-end (e2e) test scenarios. Earlier approaches in this field worked on mapping applications to page objects and creating test cases for them [52]. However, recent studies use LLMs for automated test generation [53], [54], [55]. For example, a recent prominent study [55] uses prompting patterns to track executed actions, directing the LLM to select actions based on the high-level task in the application. These actions are then formulated into test cases for the overarching task.

## VIII. CONCLUSION

Automated E2E test case generation remains challenging. In this paper, we formally defined the problem, introduced a novel methodology (AUTOE2E) for feature-driven test generation, and developed E2EBENCH, a benchmark for automated evaluation. AUTOE2E achieves **79%** feature coverage, outperforming baselines by **558%**. Future work includes enhancing AUTOE2E’s performance and exploring assertion generation for more comprehensive test coverage.

## IX. DATA AVAILABILITY

We have made AUTOE2E and E2EBENCH publicly available [56] to facilitate reproducibility of our results. Detailed instructions for replicating our experimental setup are also provided.

## REFERENCES

- [1] F. Ricca, M. Leotta, and A. Stocco, “Three open problems in the context of E2E web testing and a vision: NEONATE,” *Advances in Computers*, vol. 113, pp. 89–133, 2019.
- [2] “Selenium,” <https://www.selenium.dev>, 2024, accessed: 2024-06-28.
- [3] X. Chang, Z. Liang, Y. Zhang, L. Cui, Z. Long, G. Wu, Y. Gao, W. Chen, J. Wei, and T. Huang, “A Reinforcement Learning Approach to Generating Test Cases for Web Applications,” in *2023 IEEE/ACM International Conference on Automation of Software Test (AST)*, 2023, pp. 13–23.
- [4] S. Yu, C. Fang, X. Li, Y. Ling, Z. Chen, and Z. Su, “Effective, platform-independent gui testing via image embedding and reinforcement learning,” *ACM Transactions on Software Engineering and Methodology*.
- [5] R. K. Yandrapally and A. Mesbah, “Fragment-Based Test Generation for Web Apps,” *IEEE Transactions on Software Engineering*, vol. 49, no. 3, pp. 1086–1101, 2023.
- [6] M. Biagiola, A. Stocco, F. Ricca, and P. Tonella, “Dependency-Aware Web Test Generation,” in *2020 IEEE 13th International Conference on Software Testing, Validation and Verification (ICST)*, 2020, pp. 175–185.

- [7] Z. Xie, Y. Chen, C. Zhi, S. Deng, and J. Yin, "Chatunitest: a chatgpt-based automated unit test generation tool," *arXiv preprint arXiv:2305.04764*, 2023.
- [8] S. Kang, J. Yoon, and S. Yoo, "Large language models are few-shot testers: Exploring llm-based general bug reproduction," in *Proceedings of the 45th International Conference on Software Engineering*. IEEE Press, 2023, p. 2312–2323.
- [9] C. Lemieux, J. P. Inala, S. K. Lahiri, and S. Sen, "CodaMosa: Escaping Coverage Plateaus in Test Generation with Pre-Trained Large Language Models," in *Proceedings of the 45th International Conference on Software Engineering*, ser. ICSE '23. IEEE Press, 2023, p. 919–931.
- [10] N. Nashid, M. Sintaha, and A. Mesbah, "Retrieval-Based Prompt Selection for Code-Related Few-Shot Learning," in *Proceedings of the 45th International Conference on Software Engineering*, ser. ICSE '23. IEEE Press, 2023, p. 2450–2462.
- [11] M. Schäfer, S. Nadi, A. Eghbali, and F. Tip, "Adaptive test generation using a large language model," *arXiv preprint arXiv:2302.06527*, 2023.
- [12] Z. Liu, C. Chen, J. Wang, M. Chen, B. Wu, X. Che, D. Wang, and Q. Wang, "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*, 2024, pp. 1–13.
- [13] Z. Liu, C. Chen, J. Wang, M. Chen, B. Wu, Z. Tian, Y. Huang, J. Hu, and Q. Wang, "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*, 2024, pp. 1–12.
- [14] P. Alian, N. Nashid, M. Shahbandeh, and A. Mesbah, "Semantic constraint inference for web form test generation," in *Proceedings of the 33rd ACM SIGSOFT International Symposium on Software Testing and Analysis*, ser. ISSTA 2024. New York, NY, USA: Association for Computing Machinery, 2024, p. 932–944. [Online]. Available: <https://doi.org/10.1145/3650212.3680332>
- [15] 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.
- [16] "The claude 3 model family: Opus, sonnet, haiku," [https://www.cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbcb618857627/Model\\_Card\\_Claude\\_3.pdf](https://www.cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbcb618857627/Model_Card_Claude_3.pdf), 2024, accessed: 2024-07-26.
- [17] "Amazon web application," <https://www.amazon.com/>, 2024, accessed: 2024-07-10.
- [18] "Cypress," <https://www.cypress.io/>, 2024, accessed: 2024-06-28.
- [19] "Playwright," <https://playwright.dev/>, 2024, accessed: 2024-06-28.
- [20] A. Mesbah, A. Van Deursen, and S. Lenselink, "Crawling ajax-based web applications through dynamic analysis of user interface state changes," *ACM Transactions on the Web (TWEB)*, vol. 6, no. 1, pp. 1–30, 2012.
- [21] M. Biagiola, A. Stocco, F. Ricca, and P. Tonella, "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*, 2019, pp. 142–153.
- [22] H. Chase, "Langchain," October 2022. [Online]. Available: <https://github.com/langchain-ai/langchain>
- [23] A. Neelakantan, T. Xu, R. Puri, A. Radford, J. M. Han, J. Tworek, Q. Yuan, N. Tezak, J. W. Kim, C. Hallacy *et al.*, "Text and code embeddings by contrastive pre-training," *arXiv preprint arXiv:2201.10005*, 2022.
- [24] R. K. Yandrapally and A. Mesbah, "Fragment-based test generation for web apps," *IEEE Transactions on Software Engineering*, vol. 49, no. 3, pp. 1086–1101, 2022.
- [25] R. Yandrapally, S. Sinha, R. Tzoref-Brill, and A. Mesbah, "Carving ui tests to generate api tests and api specification," in *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*. IEEE, 2023, pp. 1971–1982.
- [26] A. Stocco, M. Leotta, F. Ricca, and P. Tonella, "Clustering-aided page object generation for web testing," in *Web Engineering: 16th International Conference, ICWE 2016, Lugano, Switzerland, June 6-9, 2016. Proceedings 16*. Springer, 2016, pp. 132–151.
- [27] M. Biagiola, A. Stocco, A. Mesbah, F. Ricca, and P. Tonella, "Web test dependency detection," in *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2019, pp. 154–164.
- [28] Y. Pan, D. Kong, S. Zhou, C. Cui, Y. Leng, B. Jiang, H. Liu, Y. Shang, S. Zhou, T. Wu *et al.*, "Webcanvas: Benchmarking web agents in online environments," *arXiv preprint arXiv:2406.12373*, 2024.
- [29] A. Drouin, M. Gasse, M. Caccia, I. H. Laradji, M. Del Verme, T. Marty, D. Vazquez, N. Chapados, and A. Lacoste, "WorkArena: How capable are web agents at solving common knowledge work tasks?" in *Proceedings of the 41st International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, vol. 235. PMLR, 21–27 Jul 2024, pp. 11 642–11 662.
- [30] AutoGPT Team, "AutoGPT," <https://github.com/Significant-Gravitas/AutoGPT>, 2024, accessed: 2024-07-30.
- [31] OpenDevin Team, "OpenDevin: An Open Platform for AI Software Developers as Generalist Agents," <https://github.com/OpenDevin/OpenDevin>, 2024, accessed: 2024-07-30.
- [32] OpenAI, "Hello gpt-4o," <https://openai.com/index/hello-gpt-4o/>, 2024, accessed: 2024-12-19.
- [33] Anthropic, "Claude 3.5 sonnet," <https://www.anthropic.com/news/claude-3-5-sonnet>, 2024, accessed: 2024-12-09.
- [34] 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*. PMLR, 2022, pp. 9466–9482.
- [35] Y. Li and O. Riva, "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*, 2021, pp. 1420–1430.
- [36] S. Mazumder and O. Riva, "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)*, 2021.
- [37] S. Jia, J. Kiros, and J. Ba, "Dom-q-net: Grounded rl on structured language," *arXiv preprint arXiv:1902.07257*, 2019.
- [38] N. Xu, S. Masling, M. Du, G. Campagna, L. Heck, J. Landay, and M. Lam, "Grounding open-domain instructions to automate web support tasks," in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, K. Toutanova, A. Rumshisky, L. Zettlemoyer, D. Hakkani-Tur, I. Beltagy, S. Bethard, R. Cotterell, T. Chakraborty, and Y. Zhou, Eds. Association for Computational Linguistics, 2021, pp. 1022–1032.
- [39] I. Gur, U. Rueckert, A. Faust, and D. Hakkani-Tur, "Learning to navigate the web," *arXiv preprint arXiv:1812.09195*, 2018.
- [40] E. Z. Liu, K. Guu, P. Pasupat, T. Shi, and P. Liang, "Reinforcement learning on web interfaces using workflow-guided exploration," *arXiv preprint arXiv:1802.08802*, 2018.
- [41] Y. Zheng, Y. Liu, X. Xie, Y. Liu, L. Ma, J. Hao, and Y. Liu, "Automatic web testing using curiosity-driven reinforcement learning," in *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*. IEEE, 2021, pp. 423–435.
- [42] L. Wang, C. Ma, X. Feng, Z. Zhang, H. Yang, J. Zhang, Z. Chen, J. Tang, X. Chen, Y. Lin *et al.*, "A survey on large language model based autonomous agents," *arXiv preprint arXiv:2308.11432*, 2023.
- [43] 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.
- [44] H. Furuta, O. Nachum, K.-H. Lee, Y. Matsuo, S. S. Gu, and I. Gur, "Multimodal web navigation with instruction-finetuned foundation models," *arXiv preprint arXiv:2305.11854*, 2023.
- [45] P. Sodhi, S. Branavan, and R. McDonald, "Heap: Hierarchical policies for web actions using llms," *arXiv preprint arXiv:2310.03720*, 2023.
- [46] M. Schäfer, S. Nadi, A. Eghbali, and F. Tip, "An empirical evaluation of using large language models for automated unit test generation," *IEEE Transactions on Software Engineering*, vol. 50, no. 1, pp. 85–105, 2024.
- [47] Y. Chen, Z. Hu, C. Zhi, J. Han, S. Deng, and J. Yin, "Chatunitest: A framework for llm-based test generation," in *Companion Proceedings of the 32nd ACM International Conference on the Foundations of Software Engineering*, 2024, pp. 572–576.
- [48] Z. Yuan, M. Liu, S. Ding, K. Wang, Y. Chen, X. Peng, and Y. Lou, "Evaluating and improving chatgpt for unit test generation," *Proc. ACM Softw. Eng.*, vol. 1, no. FSE, jul 2024. [Online]. Available: <https://doi.org/10.1145/3660783>
- [49] M. Taeb, A. Swearngin, E. Schoop, R. Cheng, Y. Jiang, and J. Nichols, "Axxnav: Replaying accessibility tests from natural language," in *Proceedings of the CHI Conference on Human Factors*

- in *Computing Systems*, ser. CHI '24. New York, NY, USA: Association for Computing Machinery, 2024. [Online]. Available: <https://doi.org/10.1145/3613904.3642777>
- [50] D. Zimmermann and A. Koziol, "Automating gui-based software testing with gpt-3," in *2023 IEEE International Conference on Software Testing, Verification and Validation Workshops (ICSTW)*, 2023, pp. 62–65.
  - [51] Z. Liu, C. Chen, J. Wang, X. Che, Y. Huang, J. Hu, and Q. Wang, "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)*, 2023, pp. 1355–1367.
  - [52] A. Stocco, M. Leotta, F. Ricca, and P. Tonella, "Apogen: automatic page object generator for web testing. *softw. qual. j.* 25 (3), 1007–1039 (2016)."
  - [53] Z. Wang, W. Wang, Z. Li, L. Wang, C. Yi, X. Xu, L. Cao, H. Su, S. Chen, and J. Zhou, "Xuat-copilot: Multi-agent collaborative system for automated user acceptance testing with large language model," *arXiv preprint arXiv:2401.02705*, 2024.
  - [54] Z. Liu, C. Chen, J. Wang, M. Chen, B. Wu, X. Che, D. Wang, and Q. Wang, "Chatting with gpt-3 for zero-shot human-like mobile automated gui testing. *arxiv* 2023," *arXiv preprint arXiv:2305.09434*.
  - [55] —, "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*, ser. ICSE '24. Association for Computing Machinery, 2024.
  - [56] "Autoe2e and e2ebench," <https://github.com/parsaalian/autoe2e>, 2024.