

# Model-based Automated Testing of JavaScript Web Applications via Longer Test Sequences

Pengfei Gao and Fu Song

School of Information Science and Technology  
ShanghaiTech University

Taolue Chen

School of computer Science and Engineering  
Nanyang Technological University

Yao Zeng

School of Computer Science and Software Engineering  
East China Normal University

Ting Su

School of computer Science and Engineering  
Nanyang Technological University

**Abstract**—JavaScript has become one of the most widely used languages for Web development. However, it is challenging to ensure the correctness and reliability of Web applications written in JavaScript, due to their dynamic and event-driven features. A variety of dynamic analysis techniques for JavaScript Web applications have been proposed, but they are limited in either coverage or scalability. In this paper, we propose a model-based automated approach to achieve high code coverage in a reasonable amount of time via testing with longer event sequences. We implement our approach as the tool LJS, and perform extensive experiments on 21 publicly available benchmarks (18,559 lines of code in total). On average, LJS achieves 86.4% line coverage in 10 minutes, which is 5.4% higher than that of JSDEP, a breadth-first search based automated testing tool enriched with partial order reduction. In particular, on large applications, the coverage of LJS is 11-18% higher than that of JSDEP. Our empirical finding supports that longer test sequences can achieve higher code coverage in JavaScript testing.

## I. INTRODUCTION

JavaScript is a highly dynamic programming language with first-class functions and “no crash” philosophy, which allows developers to write code without type annotations, and to generate and load code at runtime. Partially due to these programming flexibilities, Web applications based on JavaScript are gaining increasing popularity. These features are, however, double-edged sword, making these Web applications prone to errors and intractable to static analysis.

Dynamic analysis has proven to be an effective way to test JavaScript Web applications [1], [2], [3], [4], [5], [6], [7], [8], [9], [10], [11], [12], [13], [14], [15], [16], [17], [18], [19], [20], [21]. Since it requires testcases to explore the state space of the application, various approaches for automated testcase generation have been developed in literature, which can generate event sequences and/or associated input data of events. The event sequences concern the order in which event handlers are executed (e.g., the order of clicking buttons), while the input data concerns the choice of values (e.g., strings, numbers and objects). The generation of both event sequences and input data is important to achieve a high code coverage, and has been extensively studied.

In general, event sequences are generated by randomly selecting event handlers with heuristic search strategies [1],

[2], [3], [10], [13]. These approaches are able to analyze large real-life applications, but are usually left with a low code coverage. One possible reason, as mentioned before [12], [22], is that the event sequences are insufficiently long to explore parts of the code which may trigger the error. For example, in an experiment of [12], the uncovered code of the benchmark *tetris* is mainly due to the function *gameOver* which will only be invoked after a long event sequence. For traditional white-box testing and GUI testing, it has been shown that increasing the length of test cases could improve coverage and failure detection [23], [24], [25], [26], [27]. However, this has not been fully exploited in testing JavaScript Web applications. One of the reasons is that existing approaches usually generate event sequences from scratch by iteratively appending events to the constructed sequences up to a maximum bound, and the number of event sequences may blow up exponentially in terms of this bound. Therefore, for efficiency consideration, the maximum bound often have to be small (for instance, less than 6 [12]). To mitigate these issues, pruning techniques (e.g. mutation testing [7], [15] and partial order reduction [22]) were proposed to remove redundant event sequence, which allow to explore *limited* longer event sequences in a reasonable time. On the other hand, the input data is generated by either randomly choosing values with lightweight heuristic strategies [2], [3], [5], or using heavyweight techniques (e.g., symbolic/concolic execution) [1], [6], [11], [12], [16], [17], [19]. These works either consider unit testing or usually simply reuse the aforementioned methods to generate event sequences.

In this work, we focus on the event sequence generation. In particular, we propose a novel model-based automated testing approach to achieve a high code coverage in a reasonable time by generating and executing long event sequences. Our approach mainly consists of two key components: model constructor and event sequence generator. The model constructor iteratively queries an execution engine to generate a finite-state machine (FSM) model. It explores the state space using long event sequences and in a way to avoid prefix event subsequences re-executing and backtracking. To improve scalability, we propose a state abstraction approach, as well as a weighted

event selection strategy, to construct small-sized FSM models. The event sequence generator creates long event sequences by randomly traversing the FSM model from the initial state. We implement our approach in a tool **Longer JavaScript (LJS)**. To compare with other methods, we also implemented a random event selection strategy, and event sequence generator from JSDEP [22] based on the FSM model.

In summary, the contributions of this paper are

- A new model-based automated testing approach for client-side JavaScript Web applications via longer event sequences;
- An implementation of our approach as a tool LJS; and
- An evaluation of 21 benchmarks from JSDEP [22] including 18,559 lines of code in total.

LJS is available from <https://github.com/JavaScriptTesting/LJS>. On average, it achieves 86.4% line coverage in 10 minutes and is 5.4% higher than that of the tool JSDEP [22]. On large applications, the coverage of LJS is 11-18% higher than that of JSDEP. We have found that long event sequences can indeed improve the coverage with respect to the application under test. This paper then provides concrete, empirically validated approaches to generate long event sequences.

*Structure.* Section II introduces preliminaries, a running example and an overview of LJS. Section III describes the methodology of LJS. Section IV presents experimental results. Section V discusses related work. We conclude in Section VI.

## II. PRELIMINARIES AND OVERVIEW

### A. Preliminaries

1) *JavaScript Web Application:* A client-side Web application, which is the main focus of the paper, consists of HTML/Script files which are executed by a Web browser. When a browser loads a Web page, it parses the HTML/Script files, represents them as a Document Object Model (DOM) tree, and then executes the top-level script code. Each node in the DOM tree represents an object on the Web page and may also be associated with a set of events. Each event may have some event handlers (e.g., callback functions such as *onload* and *onclick*) which are either statically registered inside the HTML file or dynamically registered by executing functions. When an event occurs (e.g., a button is clicked), the corresponding event handlers are executed sequentially. Although the browser ensures that each callback function is executed atomically, the execution of the entire Web application exhibits nondeterminism due to the interleaving of executions of multiple callback functions.

2) *DOM Event Dependency:* Given a JavaScript Web application, let  $R_c$  and  $R_d$  respectively denote the *control* and *data* dependency relation over functions of the application. Formally, for each pair of functions (usually event handlers)  $m_1$  and  $m_2$ ,  $(m_1, m_2) \in R_c$  (resp.  $(m_1, m_2) \in R_d$ ) if there are two statements  $st_1$  in  $m_1$  and  $st_2$  in  $m_2$  that  $st_1$  affects the control (resp. data) flow of  $st_2$ . Given two DOM events  $e_1$  and  $e_2$ ,  $e_2$  is *dependent on*  $e_1$ , denoted by  $e_1 \rightarrow e_2$ , if (1) there are event handlers  $m_1$  and  $m_2$  of  $e_1$  and  $e_2$  respectively

```

1 <!DOCTYPE html>
2 <html>
3 <head>
4   <p> Example with 3 checkboxes and 1 button </p>
5 </head>
6 <body>
7   <div id="checkboxes">
8     <input id="A" type="checkbox" onclick="FA(this)"> A
9     <input id="B" type="checkbox" onclick="FB(this)"> B
10    <input id="C" type="checkbox" onclick="FC(this)"> C
11  </div>
12  <button id="Submit" type="button" > Submit </button>
13 </script>
14   var count = 0;
15   function FA(node) {
16     if (node.checked == false) count = count - 1;
17     else count = count + 1;
18     CheckedEnough(); }
19   function FB(node) {
20     if (node.checked == false) count = count - 1;
21     else count = count + 1;
22     CheckedEnough(); }
23   function FC(node){
24     if (node.checked == false) count = count - 1;
25     else count = count + 1;
26     CheckedEnough(); }
27   function CheckedEnough() {
28     var b = document.getElementById("Submit");
29     if (count >= 3) b.onclick = FSubmit;
30     else b.onclick = null; }
31   function FSubmit() {alert("Submit successfully");}
32 </script>
33 </body>
34 </html>

```

Fig. 1: Example HTML page and associated JavaScript code.

such that  $(m_1, m_2) \in (R_c \cup R_d)^*$ , or (2) the execution of  $m_1$  registers, removes, or modifies  $m_2$ . Two event sequences  $\rho_1$  and  $\rho_2$  are *equivalent* if  $\rho_1$  can be transformed from  $\rho_2$  by repeatedly swapping adjacent and independent events of  $\rho_2$ .

3) *Finite State Machine:* A (nondeterministic) Finite State Machine (FSM) is a tuple  $M = (S, I, \delta, s_0)$ , where  $S$  is a finite set of states with  $s_0 \in S$  as the initial state,  $I$  is a finite input alphabet,  $\delta \subseteq S \times I \times S$  is the transition relation. A transition  $(s_1, e, s_2) \in \delta$  denotes that after reading the input symbol  $e$  at the state  $s_1$ , the FSM  $M$  can move from the state  $s_1$  to the state  $s_2$ . We denote by  $\text{supp}(s)$  the set  $\{(e, s') \in I \times S \mid (s, e, s') \in \delta\}$ . Given a word  $e_1 \cdots e_n \in I^*$ , a *run* of  $M$  on  $e_1 \cdots e_n$  is a sequence of states  $s_0 s_1 \cdots s_n$  such that for every  $1 \leq i \leq n$ ,  $(s_{i-1}, e_i, s_i) \in \delta$ . We denote by  $\epsilon$  the empty word. Note that, in this work, we will use an FSM to represent the behaviors of a JavaScript Web application. Because of this, we do not introduce the final states as in classic FSMs.

### B. Running Example

Consider the running example in Figure 1, where the HTML code defines the DOM elements of three checkboxes and one button. The JavaScript code defines a global variable `count` and five functions manipulating `count` and the DOM elements. Initially, the three checkboxes, named A–C, are unchecked; the button, named `Submit`, does not have any *onclick* event handler. The *onclick* event handlers of the three checkboxes are the functions FA–FC respectively. For each  $X \in \{A, B, C\}$ , when the checkbox  $X$  is clicked (i.e., the corresponding event occurs), its state (checked/unchecked) is switched, and then the *onclick* event handler, namely the function  $FX$ , is executed.  $FX$  first determines whether the state

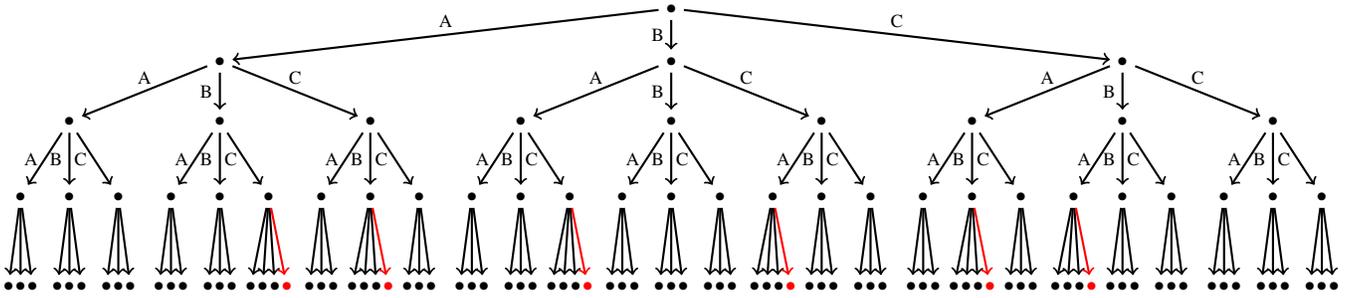


Fig. 2: The search tree of our running example, where the missed labels of edges in black color are respectively A, B, C, and the missed labels of edges in red color are Submit.

of X is checked or not. If it is checked, then the global variable count increases by one, otherwise count decreases by one. Finally, FX invokes the function CheckedEnough to verify whether the value of count is no less than 3, i.e., there are at least three checkboxes in the checked state. If  $\text{count} \geq 3$ , then the FSubmit function is registered to the *onclick* event handler of the button Submit; otherwise (i.e.,  $\text{count} < 3$ ), the *onclick* event handler of the button Submit is removed. If the button Submit is clicked when FSubmit serves the *onclick* event handler of Submit, FSubmit gets executed and prints a message to the console.

In this example, for  $X, Y \in \{A, B, C\}$ , the *onclick* event of X is dependent upon the *onclick* event of Y, and the *onclick* event of Submit is dependent upon the *onclick* event of X.

### C. Limitations of Existing Approaches

We now demonstrate why long event sequences are important for automated testing of JavaScript Web applications. Tools like ARTEMIS [2] and JSDEP [22] generate event sequences by systematically triggering various DOM events up to a fixed depth. After loading the Web page, these tools start by exploring all the available events at an initial state. If a new state is reached by executing a sequence of events, then all the available events at the new state are appended to the end of the event sequence. The procedure is repeated until time out or a fixed depth is reached. The search tree of our running example up to depth four is depicted in Figure 2, where the unshown labels of edges in black are respectively A, B, C, and those in red are Submit. Each edge labeled by  $X \in \{A, B, C, \text{Submit}\}$  denotes the execution of the *onclick* event handler of X, and each node denotes a state. For each event sequence  $\rho$  of length three, if  $\rho$  is a permutation of A;B;C, then there are four available events A, B, C, Submit after  $\rho$ , otherwise there are three available events A, B, C.

A naïve algorithm (like the default algorithm in ARTEMIS) would inefficiently explore the event space, i.e., the full tree in Figure 2, and may generate  $6 + \sum_{i=1}^4 3^i = 126$  event sequences. However, many of them are redundant. For instance, the sequences A;B;C;Submit and B;A;C;Submit actually address the same part of the code, hence one of them is unnecessary for testing purposes. To remedy this issue, JSDEP implemented a partial-order reduction in ARTEMIS

which prunes redundant event sequences by leveraging DOM event dependencies.

To cover all code of the running example, each event handler of all three checkboxes has to be executed at least two times (examining checked/unchecked states). Therefore, a sequence of length seven (e.g., A;B;C;Submit;A;B;C) is sufficient to fully cover all the code. However, if one sets the depth bound of the test sequence to be seven for the full code coverage, the default search algorithm in ARTEMIS and JSDEP may explore at least  $\sum_{i=1}^7 3^i = 3279$  event sequences. Notice that both ARTEMIS and JSDEP may re-execute previously executed test sequences in order to explore further the event space, which is time-consuming. Partially because of this, within a time limit these tools often generate and execute only short event sequences. Similar to classic program analysis, it is not hard to envision that short event sequences would hamper the coverage of code. Indeed in our running example, covering the function Submit requires test sequences with length at least 4. Unfortunately, existing approaches suffer from the “test sequence explosion” problem when increasing the depth bound of testing sequences. In this work, we propose a model-based, automated testing approach for JavaScript Web applications, aiming to generate long event sequences to improve the code coverage, but do so in a clever way to mitigate the issue of exponential blowup.

### D. Overview of Our Approach

Figure 3 presents an overview of our approach implemented in LJS, which consists of four components: static analysis, execution engine, model construction and testcase generation. Given the HTML/JavaScript source file(s) of a (client-side JavaScript Web) application, a length bound of event sequences (Max. Length), and a bound of the number of event sequences to be generated (#Testcases) as inputs, LJS outputs a (line) coverage report. Internally, LJS goes through the following steps: (1) compute the DOM event dependencies via static analysis. (2) construct an FSM model of the application with a variant of depth-first search by leveraging an execution engine and the DOM event dependencies. The FSM model is used to generate long event sequences. (3) all the generated event sequences are iteratively executed on the execution engine and output the coverage report. We now elaborate these steps in more detail.

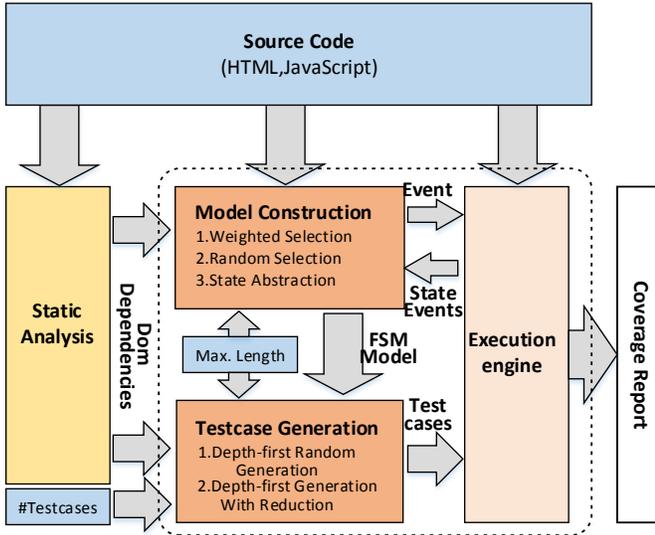


Fig. 3: Framework overview of LJS.

**Static Analysis.** Given the HTML/JavaScript source file(s) of an application, this component computes the DOM event dependencies. In our implementation, we leverage JSDEP [22] to compute DOM dependency. It first constructs a control flow graph (CFG) of the JavaScript code, and then traverses the CFG and encodes the control and data flows (i.e. dependency relations) in Datalog. The DOM event dependencies are finally computed via a Datalog inference engine. More details can be found in [22]. The analysis can handle dynamic registration, triggering, and removal of event handlers, but not other dynamic nature like dynamic code injection and obfuscation.

**Execution Engine.** Execution engine is used for FSM model construction and event sequence execution. It loads and parses the source file(s), and then executes the top-level JavaScript code. In particular, for FSM model construction, it interacts with the model constructor by iteratively receiving input events and outputting an (abstract) successor state and a set of available events at the successor after executing the input event. For event sequence execution, it receives a set of event sequences, executes them one by one and finally outputs a coverage report. For these purpose, we implement an execution engine based on ARTEMIS [2] with its features such as event-driven execution model, interaction with DOM of web pages, dynamic event handlers detection.

**Model Construction.** The model constructor interacts with the execution engine by iteratively making queries to generate an FSM model up to the given length bound (i.e., Max. Length). The FSM model is intended to represent behaviors of the application. A state of the FSM model denotes an (abstract) state of the application, and a transition  $(s, e, s')$  denotes that after executing the event handler  $e$  at the state  $s$ , the application enters the state  $s'$ . The model constructor starts with an FSM containing only one initial state, explores new state  $s'$  by selecting one event  $e$  available at the current state  $s$ , adds  $(s, e, s')$  into the FSM model, and continues exploring the state  $s'$ . LJS allows to restart the exploration from the initial state

### Algorithm 1 Model Construction

**Input:** An application  $P$  and a Max. bound  $d$

**Output:** An FSM model  $M = (S, I, \delta, s_0)$

```

1:  $i := 0; I := \emptyset; \delta := \emptyset;$ 
2:  $cur := \text{GetInitPage}(P);$ 
3:  $s_0 := \text{GetState}(cur);$ 
4:  $S := \{s_0\};$ 
5: while  $i < d$  do
6:    $e := \text{GetEvent}(cur);$ 
7:    $suc := \text{GetNextPage}(e);$ 
8:    $s' := \text{GetState}(suc);$ 
9:    $\delta := \delta \cup \{(s_0, e, s')\}; I := I \cup \{e\};$ 
10:   $S := S \cup \{s'\}; s_0 := s';$ 
11:   $cur := suc; i := i + 1;$ 
12: return  $(S, I, \delta, s_0);$ 

```

and adds new states and transitions into the FSM model, in order to make the FSM model more complete.

**Testcase Generation.** The testcase generator traverses the FSM model to generate event sequences. LJS supports two event sequence generation algorithms: (1) partial-order reduction (POR) based event sequence generation, and (2) random event sequence generation. The first algorithm traverses the FSM model from the initial state and covers all the paths up to a (usually small) bound, while redundant event sequences are pruned based on the POR from JSDEP [22]. It usually generates many short event sequences and is regarded as the baseline algorithm. The second algorithm repeatedly and randomly traverses the FSM model from the initial state to generate a small number of longer (up to a usually large bound) event sequences, and, as such, does not cover all possible paths.

To give a first impression of the performance of LJS, we ran ARTEMIS, JSDEP and LJS on the running example introduced in Section II-B. LJS reached 100% (line) coverage in 0.2s using one event sequence with length 7, ARTEMIS executed 3209 event sequences in 10min and reached 86%, and JSDEP executed 64 event sequences in 0.7s and reached 100% coverage. We also ran JSDEP and LJS on a variant of the running example which contains 10 checkboxes and the condition count  $\geq 3$  is replaced by count  $\geq 6$ . LJS reached 100% coverage in 0.4s using one event sequence with length  $(10 \times 2 + 1) = 21$ , while JSDEP executed 462 event sequences in 27.4s and reached 100% coverage. This suggests that a small number of longer event sequence could outperform a large number of shorter event sequence for applications with intensive DOM event dependency.

## III. METHODOLOGY

In this section, we present details of our model construction and testcase generation procedures.

### A. Model Construction

Algorithm 1 presents the model construction procedure, which takes an application  $P$  and a maximum bound  $d$  as inputs, and outputs an FSM model  $M = (S, I, \delta, s_0)$ . It first calls the function `GetInitPage` which loads and parses  $P$ ,

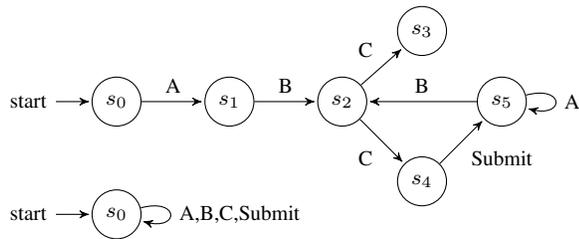


Fig. 4: The FSM models of the running example: the top-part (resp. bottom-part) FSM is constructed using the state abstraction from ARTEMIS (resp. our new state abstraction).

then executes the top-level JavaScript code, finally returns the initial Web page  $cur$  and the set of available events at this page which are dynamically detected. The initial state of the FSM model  $s_0$  is obtained by calling the function  $GetState(cur)$ . Intuitively,  $GetState$  computes a state from the source code of the Web page  $cur$  (see below). After the initialization, Algorithm 1 iteratively selects and executes events to explore the state space up to  $d$  rounds. During each iteration, it selects one available event  $e$  at the current Web page  $cur$  by using the  $GetEvent$  function (Line 6), based on some event selection strategy (see below). Then, it calls  $GetNextPage$  to get the next Web page  $suc$  by executing  $e$ . The state  $s'$  of the new Web page  $suc$  is also obtained by calling  $GetState$ . Finally, a transition  $(s, e, s')$  is added into the FSM  $M$ . Variables  $cur$  and  $i$  are then updated accordingly. As mentioned before, LJS allows to restart the exploration of the state space from the initial state, so this process may be repeated many times.

Overall, our model construction explores state space in a depth-first fashion with a large length bound (Max. Length in Figure 3), and can avoid re-executing previously executed event sequences (like ARTEMIS [2] and JSDEP [22]) or tracking state changes (like CRAWLJAX [3]), but not all the paths are explored for efficiency consideration.

1) *State Abstraction*: As mentioned, states of FSM are used to represent the states of the application and are computed from Web pages via the function  $GetState$ . State abstraction is crucial to describe application states. On the one hand, states of FSM should contain enough data to distinguish different application states, as unexplored application states may be wrongly skipped if their abstracted states were explored before, which may happen if a coarse-grained state abstraction is adopted. On the other hand, over fine-grained state abstraction may generate states which are indistinguishable wrt some coverage criteria, resulting in an explosive or even infinite state-space [28]. Therefore, the implementation of  $GetState$  requires a balance between precision and scalability.

In ARTEMIS, the state abstraction is implemented by computing the hash value of the Web page, including Web page layout and dynamically updated DOM tree, but excluding the concrete values of CSS properties and application variables, and server-side states (unless an initial server state is provided by the user). It was demonstrated that this state abstraction is effective and efficient [2].

However, our experiments find it is still too fine-grained. ARTEMIS typically assigns random values to attributes of DOM nodes which are taken into account in the state abstraction. In our running example, ARTEMIS assigns a random Boolean value to the *implicit* attribute `Value` of each checkbox when testing the running example, though the value of this attribute does not affect the code coverage. (Here “implicit” means that the attribute is not explicitly given in source code, but the object has such an attribute.) The FSM model constructed using the state abstraction from ARTEMIS is depicted in Figure 4(top-part). To prevent a prohibitively large (sometimes even infinite) state-space and improve efficiency, we use a state abstraction based on the state abstraction of ARTEMIS but discarding the random values of the implicit attributes in the DOM tree. This allows to focus on the state changes made by executing events rather than random value assignments. The FSM model constructed by our state abstraction is depicted in Figure 4(bottom-part), which is much smaller in size.

Experimental results have confirmed that our state abstraction approach significantly reduces the size of FSM models with a comparable line coverage (cf. Section IV-C).

2) *Event Selection Strategies*: It is common that several events are available at a Web page. Evidently the selection of events will affect the quality of the FSM model. In LJS, we implement two event selection strategies: (1) random event selection as the baseline algorithm and (2) weighted event selection. The former randomly chooses one of the available events, namely,  $GetEvent(cur)$  returns a random available event from the Web page  $cur$ . The latter captures the impacts of the previously executed events and DOM event dependencies. Technically we focus on:

- **Frequency of Event Execution.** In principle, all the events should be given opportunities to execute. As a result, an event which has been executed would have a lower priority to be selected in the subsequent exploration. We mention that this natural idea was already used to construct an FSM model of Android app for automated testing (e.g. [29]).
- **DOM event dependency.** Some corner-case code may be explored only by some specific event sequences due to their dependency. To expose the corner code using long event sequences, the events that *depend upon the previous selected events* deserve a higher chance to be selected.

In the weighted event selection strategy, each event  $e$  is associated with a weight, which is adjustable dynamically at runtime. The weight of  $e$  is defined as follows:

$$\text{weight}(e) = \frac{\alpha_e \times x + \beta_e \times (1-x)}{N_e + 1},$$

where  $\alpha_e$  and  $\beta_e$  are weight parameters,  $x$  is a Boolean flag determined by DOM event dependency ( $x = 1$  if  $e$  depends upon the previous selected event, 0 otherwise),  $N_e$  is number of times that  $e$  has been executed.

With the weighted event selection strategy,  $GetEvent(cur)$  randomly returns one of events which have the highest weight among all available events at the Web page  $cur$ . In our

---

**Algorithm 2** Baseline Event Sequence Generation with POR

---

An FSM  $M = (S, I, \delta, s_0)$   
**Input:** A bound  $d$  of the length of test sequences  
DOM event dependency relation  $\rightarrow$   
**Output:** A set of test sequences  $T$

```
1:  $T := \emptyset$ ;  $ss := \text{NewStack}()$ ;  
2:  $ss.\text{Push}(s_0)$ ;  
3:  $\text{EXPLORE}(ss)$ ;  
4: return  $T$ ;  
5: procedure  $\text{EXPLORE}(\text{Stack} : ss)$   
6:    $s := ss.\text{Top}()$ ;  
7:    $s.\text{SelectedEvent} := \text{null}$ ;  
8:   if  $ss.\text{Length}() \leq d$  then  
9:      $s.\text{done} := \emptyset$ ;  $s.\text{sleep} := \emptyset$ ;  
10:     $E := \{e \in I \mid \exists s'.(e, s') \in \text{supp}(s)\}$ ;  
11:    while  $\exists e \in E \setminus (s.\text{done} \cup s.\text{sleep})$  do  
12:       $s.\text{done} := s.\text{done} \cup \{e\}$ ;  
13:       $s.\text{SelectedEvent} := e$ ;  
14:      for all  $s' \in \{s' \in S \mid (e, s') \in \text{supp}(s)\}$  do  
15:         $s'.\text{sleep} := \{e' \in s.\text{sleep} \mid e \not\rightarrow e' \wedge e' \not\rightarrow e\}$ ;  
16:         $ss.\text{push}(s')$ ;  
17:         $\text{EXPLORE}(ss)$ ;  
18:         $s.\text{sleep} := s.\text{sleep} \cup \{e\}$ ;  
19:    if  $s.\text{SelectedEvent} = \text{null}$  then  
20:       $\rho := \epsilon$ ;  
21:      for all  $s' \in ss$  from bottom to top do  
22:         $\rho := \rho \cdot s'.\text{SelectedEvent}$ ;  
23:       $T := T \cup \{\rho\}$ ;  
24:     $ss.\text{Pop}()$ ;
```

---

experiments,  $\alpha_e$  and  $\beta_e$  are set to be 0.7 and 0.3 respectively, which are the best configuration after tuning.

Recall the example in Figure 1. At the initial state, all the available *onclick* events of A, B, C have the same weight 0.3 (note that the *onclick* event of Submit is not available therein), i.e., they have the same chance to be selected. Suppose the *onclick* event of A is selected, the weights of the *onclick* events of A, B, C are updated to  $\frac{0.7}{2}, \frac{0.7}{1}, \frac{0.7}{1}$  respectively. LJS then randomly chooses one of the *onclick* events of B, C. Suppose the *onclick* event of B is selected on this occasion. The weights of all the available *onclick* events of A, B, C become  $\frac{0.7}{2}, \frac{0.7}{2}, \frac{0.7}{1}$ , which implies that the *onclick* event of C will be selected at the next step. After that, the weights of the *onclick* events A, B, C, Submit are updated to  $\frac{0.7}{2}, \frac{0.7}{2}, \frac{0.7}{2}, \frac{0.7}{1}$  (note that the *onclick* event of Submit now becomes available).

### B. Testcase Generation

In this work, as mentioned in the introduction, we focus on the event sequence generation while the input data is chosen randomly. We first present the baseline algorithm for generating event sequences with partial-order reduction (POR) which is inspired by [22], [30], and then discuss how to generate long test sequences.

1) *Baseline Event Sequence Generation with POR*: Given the FSM  $M = (S, I, \delta, s_0)$ , a bound  $d$ , and the DOM event dependency relation  $\rightarrow$ . Algorithm 2 (excluding Lines 15 and 18) generates all possible event sequences with length up to  $d$  stored in  $T$ .

The procedure EXPLORE traverses the FSM  $M$  in a depth-first manner, where  $ss$  is a working stack storing the event sequence with the initial state  $s_0$  as the bottom element,  $E$  denotes the set of available events at state  $s$ ,  $s.\text{SelectedEvent}$  denotes the selected event at  $s$ , and  $s.\text{done}$  denotes the set of all previously selected events at  $s$ . The procedure EXPLORE first checks whether the bound  $d$  is reached (Line 8). If  $ss.\text{Length}() > d$ , then the while-loop is skipped. After that, if  $s.\text{SelectedEvent} = \text{null}$  (indicating that the sequence in  $ss$  has reached the maximum length  $d$ ), the event sequence stored in the stack  $ss$  is added to the set  $T$ . Otherwise, the while-loop will explore a previously unexplored event and invoke the procedure EXPLORE recursively. The while-loop terminates when all the available events at state  $s$  have been explored. Note that in this case,  $s.\text{SelectedEvent} \neq \text{null}$ , hence the sequence in  $ss$  will not be added to  $T$ .

With Lines 15 and 18 Algorithm 2 implements the partial-order reduction based on the notion of sleep-set [31]. In principle, it first classifies the event sequences into equivalence classes, and then explores one representative from each equivalence class. In detail, each event  $e$  explored at a state  $s$  is put into its sleep set  $s.\text{sleep}$  (Line 18). When another event  $e'$  is explored at  $s$ , the sleep set of  $s$  is copied into the sleep set of the next state  $s'$  (Line 15), if the DOM events  $e$  and  $e'$  are independent of each other. Later, each event at  $s$  will be skipped if it is in the sleep set of  $s$  (Line 11), because executing this event is guaranteed to reach a previously explored state.

2) *Long Event Sequence Generation*: Algorithm 3 shows the pseudocode of our random event sequence generation. Given the FSM model  $M = (S, I, \delta, s_0)$ , a maximum length bound  $d$  and a maximum bound  $m$  of the number of event sequences, Algorithm 3 randomly generates  $m$  number of event sequences with length  $d$ . Each iteration of the outer while-loop computes one event sequence with length  $d$  until the number of event sequences reaches  $m$ . In the inner while-loop, it starts from the initial state  $s_0$  and an empty sequence  $\rho = \epsilon$ . At each state  $s$ , the inner while-loop iteratively and randomly selects a pair  $(e, s')$  of event and state denoting that executing the event  $e$  at the state  $s$  moves to the state  $s'$ , then appends  $e$  to the end of previously computed sequence  $\rho$ . Algorithm 3 repeats this procedure until the length of  $\rho$  reaches the maximum bound  $d$ . At this moment, one event sequence is generated and stored into the set  $T$ . The outer while-loop enters its next iteration.

Algorithm 3 may not generate all the possible event sequences, and may generate redundant event sequences, but less often, due to large maximum bound. We remark that the POR technique cannot be integrated into Algorithm 3.

## IV. EXPERIMENTS

We have implemented our method as a software tool LJS. It exploits JSDEP [22] for computing DOM event dependency and a modified automated testing framework ARTEMIS [2] as the execution engine. For comparison purpose, we implemented LJS in such a way that individual techniques are modularized and can be enabled on demand. Thus, we were

**Algorithm 3** Long Event Sequence Generation

---

An FSM  $M = (S, I, \delta, s_0)$   
**Input:** A bound  $d$  of the length of test sequences  
A bound  $m$  of the number of test sequences  
**Output:** A set of test sequences  $T$

```

1:  $T := \emptyset$ ;
2: while  $|T| < m$  do
3:    $\rho := \epsilon$ ;
4:    $s := s_0$ ;
5:   while  $|\rho| < d$  do
6:      $(e, s) := \text{RandomlySelectOnePair}(\text{supp}(s))$ ;
7:      $\rho := \rho \cdot e$ ;
8:    $T := T \cup \{\rho\}$ ;
9: return  $T$ ;

```

---

able to compare the performance of various approaches with different configurations, in particular, (1) our state abstraction vs the state abstraction from [2], (2) random event selection vs weighted event selection, (3) baseline event sequence generation with POR (i.e., Algorithm 2) vs long event sequence generation (i.e., Algorithm 3). To demonstrate the efficiency and effectiveness of LJS, we compared LJS with JSDEP [22] on same benchmarks. (We note that in [22] JSDEP is shown to be superior to ARTEMIS [2], so a direct comparison between LJS and ARTEMIS [2] is excluded.)

The experiments are designed to answer the following research questions:

- RQ1.** How efficient and effective is LJS compared with JSDEP [22]?
- RQ2.** How effective is our coarse-grained state abstraction compared with the state abstraction from [2]?
- RQ3.** How effective is the weighted event selection strategy compared with the random event selection strategy?
- RQ4.** How effective is the long event sequence generation compared with the baseline algorithm?

### A. Evaluation Setup

To make comparison on a fair basis, we evaluated LJS on publicly available benchmarks<sup>1</sup> of JSDEP [22], which consist of 21 client-side JavaScript Web applications with 18,559 lines of code in total. Columns 1-2 of Table I show the name of the application and the number of lines of code. We ran all experiments on a server with a 64-bit Ubuntu 12.04 OS, Intel Xeon(R) E5-2603v4 CPU (1.70 GHz, 6 Cores), and 32GB RAM. To answer the research questions **RQ1-RQ4**, we conducted four case studies. The time used to compute the DOM event dependency is usually marginal and can be safely ignored, so is not counted in line with [22]. For statistics, we ran LJS on each application 5 times and simply took the average as the result. The coverage measure are the aggregation of that from model construction and testcase execution respectively which are separated in the last experiment.

Name	Loc	LJS			JSDEP		
		CRG.	sd	Tests Len.	CRG.	Tests	M.Len.
case1	59	100.0%	0	4679 99	100.0%	1409	705
case2	72	100.0%	0	4376 99	100.0%	3058	549
case3	165	100.0%	0	2739 99	100.0%	7811	575
case4	196	<b>87.0%</b>	0	2816 99	<b>77.9%</b>	8594	<b>500</b>
<b>frog</b>	567	<b>96.8%</b>	0.004	<b>15</b> 99	<b>84.6%</b>	<b>86</b>	<b>16</b>
cosmos	363	82.0%	0.060	322 99	79.5%	973	243
hanoi	246	89.0%	0	1303 99	82.5%	902	225
flipflop	525	97.0%	0	59 99	96.3%	284	71
<b>sokoban</b>	<b>3056</b>	<b>88.6%</b>	0.026	<b>58</b> 99	<b>77.6%</b>	<b>203</b>	<b>51</b>
wormy	570	45.4%	0.054	35 99	41.0%	323	18
chinabox	338	84.0%	0	7 99	82.3%	92	9
<b>3dmodel</b>	<b>5414</b>	<b>85.0%</b>	0	<b>8</b> 99	<b>71.5%</b>	<b>66</b>	<b>10</b>
<b>cubuild</b>	<b>1014</b>	<b>88.8%</b>	0.060	<b>7</b> 99	<b>82.8%</b>	<b>153</b>	<b>17</b>
pearlski	960	55.0%	0	72 99	54.9%	214	52
speedyeater	784	89.8%	0.010	516 99	82.1%	1497	374
gallony	300	95.0%	0	2133 99	94.5%	1611	95
fullhouse	528	92.2%	0.012	1119 99	86.3%	889	222
<b>ball_ool</b>	<b>1745</b>	<b>92.8%</b>	0.004	<b>1</b> 99	<b>74.2%</b>	<b>18</b>	<b>4</b>
harehound	468	94.8%	0.004	522 99	94.5%	1224	116
match	369	72.8%	0.004	1063 99	73.2%	4050	845
lady	820	79.0%	0	0 99	75.7%	35	8
<b>Average</b>	883.8	<b>86.4%</b>	0.011	1040.5 99	<b>81.0%</b>	1594.9	-

TABLE I: Coverage of LJS and JSDEP in 600s.

### B. RQ1: LJS vs. JSDEP

In this study, we answer **RQ1** by performing two experiments assessing effectiveness and efficiency. For effectiveness, we compare the (line) coverage of LJS with JSDEP in 600 seconds. For efficiency, we compare the time used by LJS and JSDEP to achieve the same coverage. The experiments of LJS were performed with the following configurations: Max. Length=99, random event selection strategy, our new state abstraction, running Algorithm 1 two times, and long event sequence generation. Experiments of JSDEP were performed with the setting shown in [22].

Table I shows the results of LJS and JSDEP in 600s, where Columns 3-6 (resp. Columns 7-9) show the average, standard deviation (sd) of coverage obtained in 5 runs, number and length of event sequences after running LJS (resp. JSDEP).

Overall, we can observe an increase in the average coverage from 81% achieved by JSDEP to 86.4% achieved by LJS. (We remark that the average coverage was increased from 67% achieved by ARTEMIS to 80% achieved by JSDEP in 600s [22]. In our experiment, JSDEP performed slightly better.) Perhaps more importantly, on large applications such as *frog*, *sokoban*, *3dmodel*, *cubuild* and *ball\_pool*, the coverage of LJS is 11-18% higher than that of JSDEP. We can also observe that a small amount of long event sequences could outperform a large amount of short event sequences. However, it should be emphasized that long event sequences, but of low quality, may not improve the coverage. This has been demonstrated by the results of *case4*, *speedyeater* and *fullhouse*.

Figure 5 shows the coverage that LJS and JSDEP achieved by running on the top 6 largest applications in time ranging from 60s to 600s with step size 60s. (Note that the application *lady* is excluded in this experiments because its model construction takes more than 600s; cf. Table I.) The X-axis is the execution time budget whereas Y-axis is the achieved

<sup>1</sup><https://github.com/ChunghaSung/JSdep>.

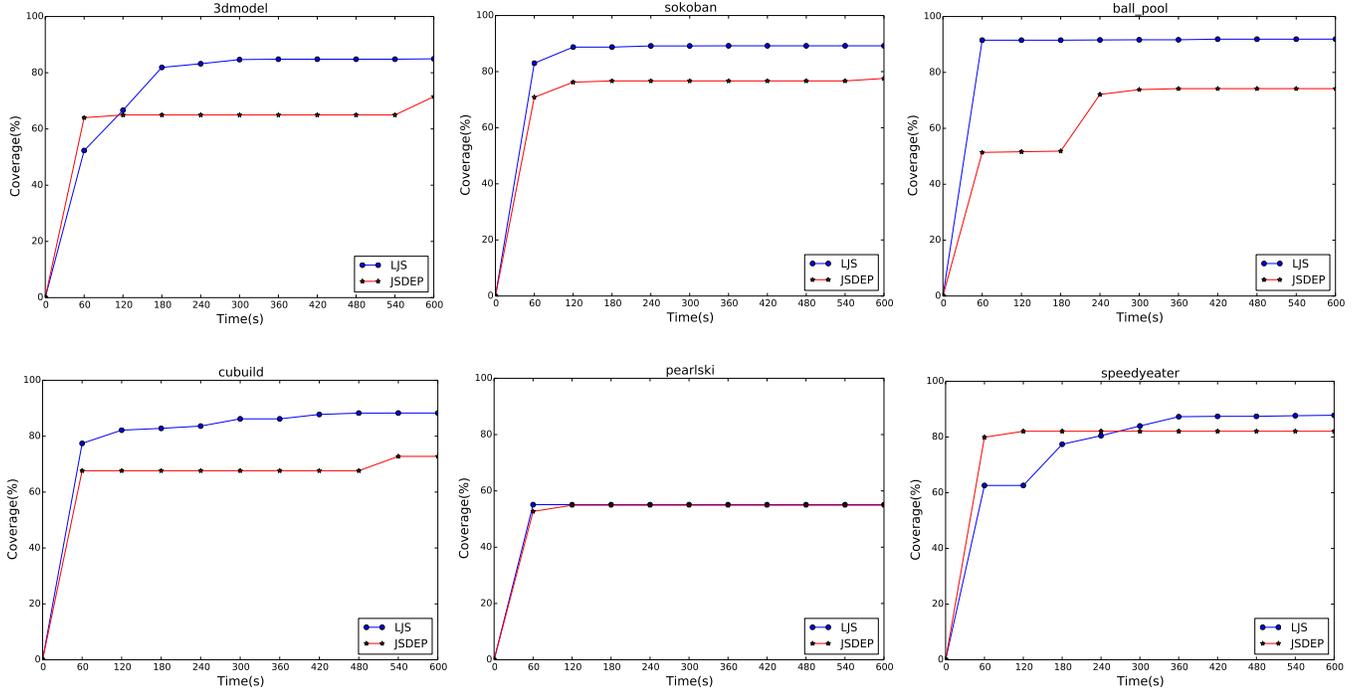


Fig. 5: Coverage of LJS and JSDEP, as a function of the execution time.

Name	State abstraction from [2]				Our state abstraction			
	CRG.	S	δ	Time(s)	CRG.	S	δ	Time(s)
case1	100.0%	1.0	2.0	0.5	100.0%	1.0	2.0	0.0
case2	97.6%	1.0	4.0	0.7	95.2%	1.0	4.0	0.9
case3	100.0%	1.0	6.0	1.2	99.0%	1.0	6.0	1.4
case4	87.0%	1.0	8.0	1.3	87.0%	1.0	8.0	1.5
<b>frog</b>	95.2%	<b>199.0</b>	<b>198.0</b>	87.5	94.8%	<b>24.2</b>	<b>100.8</b>	105.9
<b>cosmos</b>	79.0%	<b>125.0</b>	<b>196.0</b>	9.0	78.8%	<b>65.6</b>	<b>142.2</b>	10.1
hanoi	89.0%	104.8	186.0	2.0	89.0%	99.0	185.4	2.3
flipflop	97.4%	25.8	100.4	21.6	97.0%	27.2	110.8	18.5
<b>sokoban</b>	88.4%	<b>68.6</b>	<b>191.4</b>	23.6	88.6%	<b>31.2</b>	<b>125.8</b>	26.4
wormy	42.2%	189.8	198.0	37.9	40.8%	132.8	185.4	38.6
chinabox	84.0%	67.6	156.4	136.8	84.0%	69.8	161.4	151.2
<b>3dmodel</b>	81.6%	<b>3.0</b>	<b>25.6</b>	105.5	83.2%	<b>1.0</b>	<b>6.0</b>	106.2
cubuild	85.0%	87.0	166.6	131.8	84.8%	87.4	163.2	132.5
<b>pearlski</b>	55.0%	<b>136.8</b>	<b>196.4</b>	18.2	55.0%	<b>75.6</b>	<b>176.8</b>	19.2
<b>speedyeater</b>	82.0%	<b>175.8</b>	<b>198.0</b>	6.0	81.2%	<b>4.6</b>	<b>65.0</b>	9.5
gallory	95.0%	63.8	174.8	1.4	95.0%	63.4	176.0	1.7
<b>fullhouse</b>	93.0%	<b>161.6</b>	<b>198.0</b>	2.5	93.0%	<b>28.6</b>	<b>116.0</b>	3.3
ball_pool	93.0%	40.6	142.2	354.4	93.0%	39.8	150.4	346.2
<b>harehound</b>	83.8%	<b>191.0</b>	<b>198.0</b>	5.9	84.8%	<b>16.2</b>	<b>109.4</b>	9.0
match	67.4%	<b>13.6</b>	<b>171.0</b>	3.2	69.8%	<b>4.4</b>	<b>62.8</b>	5.0
lady	79.2%	<b>173.4</b>	<b>197.8</b>	936.1	79.2%	<b>85.6</b>	<b>174.0</b>	1109.6
<b>Average</b>	84.5%	<b>87.2</b>	<b>138.8</b>	89.9	84.4%	<b>41.0</b>	<b>106.3</b>	100.0

TABLE II: Comparison of state abstraction techniques.

coverage. Overall, as the execution time budget increases, the coverage of LJS is higher than that of JSDEP. Moreover, the rate of LJS to achieve a higher coverage is, in most cases, slightly higher than that of JSDEP.

### C. RQ2: Comparison of State Abstraction Techniques

In this study, we answer **RQ2** by performing one experiment and comparing the obtained FSM model and coverage results using our coarse-grained state abstraction and the state abstraction from [2] respectively. In this experiment, LJS constructs

the FSM model by running Algorithm 1 two times using Max. Length=99 and random event selection, and generates two event sequences from the FSM model. Taking into account the consumption of time, the experiment only generates two event sequences.

Table II shows the results, where Columns 2-5 (resp. Columns 6-9) show the coverage, numbers of states and transitions of the FSM model (note that we take the average of these numbers), and execution time after running LJS with the state abstraction from [2] (resp. our new state abstraction).

Overall, the numbers of states and transitions using our state abstraction are much smaller, with a dramatic decrease in some large applications such as *sokoban*, *3dmodel*, *pearlski*, *speedyeater* and *harehound*. Meanwhile, the performance of the two state abstractions (in terms of average coverage and execution time) is comparable.

### D. RQ3: Comparison of Event Selection Strategies

In this study, we answer **RQ3** by performing one experiment and comparing the obtained FSM model and coverage results using weighted event selection strategy and random event selection strategy respectively. In this experiment, LJS constructs the FSM model by running Algorithm 1 two times using Max. Length=99 and our coarse-grained state abstraction, and generates two event sequences from the FSM model.

Table III shows the results, where Columns 2-5 (resp. Columns 6-9) show the coverage, numbers of states and transitions of the FSM model, and execution time after running LJS with random event selection (resp. weighted event selection).

Name	Random Event Selection				Weighted Event Selection			
	CRG.	S	δ	Time(s)	CRG.	S	δ	Time(s)
case1	100.0%	1.0	2.0	0.5	100.0%	1.0	2.0	0.6
case2	100.0%	1.0	4.0	0.9	97.6%	1.0	4.0	1.4
case3	98.0%	1.0	6.0	1.9	100.0%	1.0	6.0	2.2
case4	87.0%	1.0	8.0	1.4	87.0%	1.0	8.0	1.7
frog	95.2%	28.6	101.0	107.1	95.2%	52.0	119.8	208.9
cosmos	79.0%	66.6	143.6	9.5	76.8%	24.0	149.0	9.0
hanoi	89.0%	102.8	186.4	2.1	89.0%	112.0	121.4	3.5
flipflop	97.0%	26.6	108.2	20.3	97.0%	17.0	55.8	50.6
sokoban	89.6%	28.6	123.2	24.8	84.0%	9.0	36.0	25.4
wormy	40.6%	131.8	186.8	37.1	41.8%	134.0	174.0	37.8
chinabox	84.0%	60.0	151.6	159.2	84.0%	72.2	159.8	129.4
<b>3dmodel</b>	<b>82.0%</b>	<b>1.0</b>	<b>6.0</b>	<b>98.1</b>	<b>84.0%</b>	<b>1.0</b>	<b>6.0</b>	<b>89.1</b>
<b>cubuild</b>	<b>84.8%</b>	<b>92.8</b>	<b>171.2</b>	<b>136.9</b>	<b>85.6%</b>	<b>85.4</b>	<b>141.4</b>	<b>130.1</b>
pearlski	55.0%	71.2	175.0	18.9	51.0%	33.0	87.0	17.5
speedyeater	87.6%	4.8	65.2	8.9	82.2%	2.0	34.0	9.9
gallony	95.0%	61.6	173.8	1.5	92.0%	30.4	127.6	2.1
<b>fullhouse</b>	<b>93.0%</b>	<b>30.6</b>	<b>123.0</b>	<b>3.1</b>	<b>75.0%</b>	<b>5.0</b>	<b>13.0</b>	<b>4.7</b>
ball_pool	93.4%	42.2	146.2	345.1	92.8%	54.6	191.0	404.9
harehound	89.4%	11.8	99.2	9.1	88.2%	3.0	34.0	10.4
match	69.4%	3.6	56.4	4.8	69.6%	4.0	50.6	5.6
lady	79.0%	82.6	171.4	1849.2	78.0%	67.4	153.6	1881.1
<b>Average</b>	<b>85.1%</b>	<b>40.5</b>	<b>105.2</b>	<b>135.3</b>	<b>83.4%</b>	<b>33.8</b>	<b>79.7</b>	<b>144.1</b>

TABLE III: Comparison of event selection strategies.

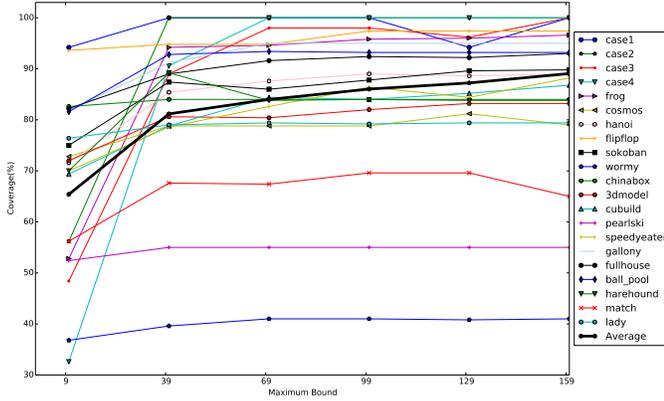


Fig. 6: Coverage of LJS, as a function of the length bound.

Overall, the numbers of states and transitions of weighted event selection are respectively less than that of random event selection. Meanwhile, the average coverage and execution time of the two selection strategies are comparable. In particular, in terms of coverage for applications such as *cubuild* the weighted event selection strategy performs better, whereas for applications such as *fullhouse*, the random event selection strategy performs better. We will discuss these findings later.

#### E. RQ4: Comparison of Event Sequence Generations

In this study, we answer **RQ4** by performing two experiments and comparing quality of event sequences generated by the long event sequence generation (Algorithm 3) and baseline event sequence generation with POR (Algorithm 2). In the first experiment, LJS constructs the FSM model by running Algorithm 1 two times using our coarse-grained state abstraction and the random event selection strategy, and generates two event sequences from the FSM model, while the maximum bounds (i.e., Max. Length) are 9, 39, 69, 99, 129 and 159. In the second experiment, we first generate an FSM model

Name	Baseline with POR				Long Sequence Generation			
	CRG.	Len.	Time(s)	Tests	CRG.	Len.	Times	Tests
case1	100.0%	16	829	65504	100.0%	99	600	4543
<b>case2</b>	<b>87.8%</b>	<b>9</b>	<b>1125</b>	<b>87040</b>	<b>100.0%</b>	<b>99</b>	<b>600</b>	<b>4244</b>
<b>case3</b>	<b>68.6%</b>	<b>7</b>	<b>931</b>	<b>54432</b>	<b>100.0%</b>	<b>99</b>	<b>600</b>	<b>2658</b>
<b>case4</b>	<b>61.8%</b>	<b>6</b>	<b>674</b>	<b>32768</b>	<b>87.0%</b>	<b>99</b>	<b>600</b>	<b>2732</b>
frog	88.6%	4	664	592	95.8%	99	600	8
cosmos	78.0%	5	1072	4415	83.8%	99	600	175
hanoi	86.8%	6	889	4276	89.0%	99	600	1365
flipflop	97.0%	6	1200	419	97.0%	99	600	66
sokoban	88.6%	5	1200	1697	89.0%	99	600	54
wormy	41.2%	13	657	317	43.2%	99	600	30
<b>chinabox</b>	<b>78.0%</b>	<b>8</b>	<b>729</b>	<b>64</b>	<b>84.0%</b>	<b>99</b>	<b>600</b>	<b>5</b>
<b>3dmodel</b>	<b>72.0%</b>	<b>4</b>	<b>1200</b>	<b>93</b>	<b>85.0%</b>	<b>99</b>	<b>600</b>	<b>10</b>
<b>cubuild</b>	<b>75.2%</b>	<b>6</b>	<b>1019</b>	<b>238</b>	<b>92.2%</b>	<b>99</b>	<b>600</b>	<b>10</b>
pearlski	52.1%	7	960	653	55.0%	99	600	85
speedyeater	82.3%	5	1200	17106	88.0%	99	600	508
gallony	94.5%	8	1200	9821	95.0%	99	600	2852
fullhouse	79.2%	8	1200	12097	79.0%	99	600	1261
ball_pool	92.1%	6	1200	22	93.0%	99	600	3
harehound	92.2%	4	1200	7949	95.0%	99	600	420
match	73.2%	5	738	13341	73.0%	99	600	1103
lady	73.2%	5	696	12	75.0%	99	600	6
<b>Average</b>	<b>79.2%</b>	<b>-</b>	<b>980</b>	<b>14897.9</b>	<b>85.7%</b>	<b>99</b>	<b>600</b>	<b>1054.2</b>

TABLE IV: Comparison of event sequence generation algorithms, where coverage (i.e. CRG) *only* includes that of event sequence execution.

using the same configuration with a fixed maximum bound 99. For each application, from the same FSM model, LJS with Algorithm 3 enabled iteratively generates and executes event sequences with maximum bound 99 within 600s time bound. Meanwhile, LJS with Algorithm 2 enabled generates and executes all event sequences with redundant sequences pruned up to some maximum bound so that the execution time exceeds 600s, moreover, the execution terminates when the execution time reaches 1200s.

Figure 6 shows the results of the first experiment, where the X-axis is the maximum bound, and Y-axis is achieved coverage with that bound, and the black bold line shows the trend of average coverage.

Overall, the average coverage increases quickly when the maximum bound increases from 9 to 39, but only slightly when it increases from 39 to 159.

Table IV shows the results of the second experiment, where Columns 2-5 (resp. Columns 6-9) show the coverage (of testcase execution *only*), maximum sequence length, execution time and number of event sequences after running LJS with Algorithm 2 (resp. Algorithm 3) enabled. Overall, the average coverage of Algorithm 3 (i.e., long event sequence generation) is 6.5% higher than that of Algorithm 2 (i.e., baseline event sequence generation with POR) with less execution time. In particular, the coverage improvement of Algorithm 3 is more prominent for applications *case2-case4*, *chinabox*, *3dmodel* and *cubuild*. These results also confirm that executing fewer long event sequences may achieve higher coverage than executing more short event sequences.

#### F. Discussions

The coverage improvement, as shown in the experiment, is not significant using the weighted event selection. There are two possible reasons. First, weighted event selection re-

lies upon DOM event dependencies which are computed by JSDEP. JSDEP uses context, path, flow and object-insensitive static analysis, hence the dependencies may be not sufficiently fine-grained. Using more accurate DOM event dependencies may improve the effectiveness of the weighted event selection. Second, the selected benchmarks might not be representative of real-world JavaScript programs that are DOM event dependency intensive. LJS currently supports off-line testing (i.e. when source code is available) when DOM event dependency is enabled. Our approach is also applicable in on-line testing if the DOM event dependency is computed dynamically, as it does not need to record and replay.

We note some limitations of the experiments. The experiments are based on a benchmark that includes only 4 large-scale Web applications (with more than 1k LOC, maximum 5k LOC). More experiments are needed to assess randomized algorithms, for instance, the statistical significance of the coverage improvement [32].

## V. RELATED WORK

We discuss the related work in the areas of model-based testing and automated JavaScript Web application testing.

### A. Model-based Automated Testing

Model-based testing (MBT) has been widely used in software testing (cf. [33], [34], [35] for surveys). Mainstream MBT techniques differ mainly in three aspects: models of the software under test, model construction, and testcase generation. Several models, such as state-based (e.g., pre-/post-condition) and transition-based (e.g., UML and I/O automata), have been proposed. The FSM model used in this work is one of the transition-based models. Model construction is one of the most important tasks in MBT. It is usually time-consuming and error-prone to manually construct models for GUI-based applications [36], [37]. Therefore, most works use static/dynamic analysis to construct models, for example, [38], [39], [40], [29], [27], [26], [41], [42] for mobile/GUI applications. However, it is rather difficult to statically construct models for JavaScript Web applications due to their dynamic characteristics [3]. Regarding works on model construction for JavaScript Web applications [43], [44], [45], [46], [3], [47]; [43], [46], [44] have to construct a model manually. [43] extracts the model via static analysis, but lacks of considering dynamic nature of JavaScript; [45], [3] construct FSM models via dynamic analysis to crawl Web applications. The main difference between our work and theirs [45], [3] is the way in which the model is constructed. Our model construction pursues larger depth without backtracking, but does not cover all possible event sequences, whereas [45], [3] cover all the possible event sequences up to a length bound with backtracking. [47] also reported to construct FSM models, but did not give detail of their algorithm, nor included JavaScript coverage. Existing testcase generation algorithms mainly focus on the systematic generation of event sequences with a rather limited length bound due to “test sequence explosion” problem. Our approach generates long event sequences, but strategically

avoid covering all possible event sequences (up to a length bound) to mitigate the exponential blowup problem. The tool STOAT [29] considered model-based testing for Android apps, but with different testcase generation strategy as ours.

### B. JavaScript Web Application Automated Testing

Web application testing has been widely studied in the past decade, differing mainly in targeted Web programming languages (e.g., PHP [48], [49], [50] and JavaScript [2], [3], [5]), and testing techniques (e.g., model-based testing [33], [34], [35], mutation testing [15], [7], [46], search-based testing [48], [51], [52], [53] and symbolic/concolic testing [1], [6], [12]; cf. [28], [21] for surveys). We mainly compare with works on automated testing of JavaScript Web Application.

Test sequences of JavaScript Web Applications consist of event sequences and input data for each event. Existing works create event sequences via exploring the state space by randomly selecting events with heuristic search strategies [1], [2], [13], [3], [10]. For instance, Kudzu [1] and CRAWLJAX [3] randomly select available events. Moreover, CRAWLJAX relies on a heuristic approach for detecting event handlers, hence may not be able to detect all of them. ARTEMIS [2] uses the heuristic strategy based on the observed read and write operations by each event handler in an attempt to exclude sequences of non-interacting event handler executions. EventBreak [13] uses the heuristic strategy based on performance costs in terms of the number of conditions in event handlers in an attempt to analyze responsiveness of the application. These approaches usually cover all the sequences up to a given, usually small length bound. In order to explore long event sequences in limited time, delta-debugging based method [54] and partial-order reduction [22] were proposed for pruning redundant event sequences. Our approach does not cover all possible event sequences, hence can create long event sequences within the time budget. Experimental results show that our approach can achieve a high line coverage than ARTEMIS even with partial-order reductions [22].

Another research line in automated testing of JavaScript Web application is to generate high quality input data of events using symbolic/concolic testing, e.g., [1], [6], [12], [16], [17], [19]. These approaches are able to achieve high coverage, but heavily rely on the underlying constraint solver. They generally do not scale well for large realistic programs, because the number of feasible execution paths of a program often increases exponentially in the length of the path. Our work focuses on the generation of long event sequences, but choose the input data randomly, which is orthogonal, and could be complementary, to the more advanced input data generation methods. A transfer technique has been proposed based on the automation engine framework SELENIUM to transfer tests from one JavaScript Web application to another [55]. This is orthogonal to this work.

## VI. CONCLUSION AND FUTURE WORK

We have proposed a model-based automated testing approach for JavaScript Web applications. Our approach distin-

guishes from others in making use of *long* event sequences in both FSM model construction and testcase generation from the FSM model. We have implemented our approach in a tool LJS and evaluated it on a collection of benchmarks. The experimental results showed that our new approach is more efficient and effective than ARTEMIS and JSDEP. Furthermore, we empirically found that longer test sequences can achieve a high line coverage.

For future work, we plan to experiment on more larger benchmarks and to have a thorough statistical analysis. In particular, the current evaluation is based on coverage, but it would also be interesting to evaluate the fault detection. Furthermore, we note that there are several ways to obtain long sequences based on FSM, for instance, by following the dependency chains or giving priority to the novel states. These strategies deserve further exploration.

## REFERENCES

- [1] P. Saxena, D. Akhawe, S. Hanna, F. Mao, S. McCamant, and D. Song, "A symbolic execution framework for JavaScript," in *Proceedings of 31st IEEE Symposium on Security and Privacy (S&P)*, 2010, pp. 513–528.
- [2] S. Artzi, J. Dolby, S. H. Jensen, A. Møller, and F. Tip, "A framework for automated testing of JavaScript web applications," in *Proceedings of the 33rd International Conference on Software Engineering (ICSE)*, 2011, pp. 571–580.
- [3] A. Mesbah, A. V. Deursen, and S. Lenseslink, "Crawling ajax-based web applications through dynamic analysis of user interface state changes," *ACM Transactions on the Web*, vol. 6, no. 1, pp. 1–30, 2012.
- [4] S. Mirshokraie and A. Mesbah, "JSART: JavaScript assertion-based regression testing," in *Proceedings of the 12th International Conference on Web Engineering (ICWE)*, 2012, pp. 238–252.
- [5] P. Heidegger and P. Thiemann, "JSConTest: Contract-driven testing and path effect inference for JavaScript," *Journal of Object Technology*, vol. 11, no. 1, pp. 1–29, 2012.
- [6] K. Sen, S. Kalasapur, T. G. Brutch, and S. Gibbs, "Jalangi: a selective record-replay and dynamic analysis framework for JavaScript," in *Proceedings of the Joint Meeting of the European Software Engineering Conference and the ACM SIGSOFT Symposium on the Foundations of Software Engineering (ESEC/FSE)*, 2013, pp. 488–498.
- [7] S. Mirshokraie, A. Mesbah, and K. Pattabiraman, "Efficient JavaScript mutation testing," in *Proceedings of the 6th IEEE International Conference on Software Testing, Verification and Validation*, 2013, pp. 74–83.
- [8] A. M. Fard and A. Mesbah, "JSNOSE: detecting JavaScript code smells," in *Proceedings of the 13th IEEE International Working Conference on Source Code Analysis and Manipulation (SCAM)*, 2013, pp. 116–125.
- [9] S. Mirshokraie, A. Mesbah, and K. Pattabiraman, "PYTHIA: generating test cases with oracles for JavaScript applications," in *Proceedings of the 28th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 2013, pp. 610–615.
- [10] S. Mirshokraie, "Effective test generation and adequacy assessment for JavaScript-based web applications," in *Proceedings of International Symposium on Software Testing and Analysis (ISSTA)*, 2014, pp. 453–456.
- [11] A. M. Fard, M. MirzaAghaei, and A. Mesbah, "Leveraging existing tests in automated test generation for web applications," in *Proceedings of the ACM/IEEE International Conference on Automated Software Engineering (ASE)*, 2014, pp. 67–78.
- [12] G. Li, E. Andreasen, and I. Ghosh, "SymJS: automatic symbolic testing of JavaScript web applications," in *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering*, 2014, pp. 449–459.
- [13] M. Pradel, P. Schuh, G. Necula, and K. Sen, "EventBreak: analyzing the responsiveness of user interfaces through performance-guided test generation," in *Proceedings of the 2014 ACM International Conference on Object Oriented Programming Systems Languages & Applications (OOPSLA)*, 2014, pp. 33–47.
- [14] A. M. Fard, A. Mesbah, and E. Wohlstader, "Generating fixtures for JavaScript unit testing," in *Proceedings of the 30th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 2015, pp. 190–200.
- [15] S. Mirshokraie, A. Mesbah, and K. Pattabiraman, "Guided mutation testing for JavaScript web applications," *IEEE Trans. Software Eng.*, vol. 41, no. 5, pp. 429–444, 2015.
- [16] K. Sen, G. C. Necula, L. Gong, and W. Choi, "MultiSE: multi-path symbolic execution using value summaries," in *Proceedings of the 10th Joint Meeting on Foundations of Software Engineering (ESEC/FSE)*, 2015, pp. 842–853.
- [17] G. L. Hideo Tanida, Tadahiro Uehara and I. Ghosh, "Automated unit testing of JavaScript code through symbolic executor SymJS," *International Journal on Advances in Software*, vol. 8, no. 1-2, pp. 146–155, 2015.
- [18] S. Mirshokraie, A. Mesbah, and K. Pattabiraman, "Atrina: Inferring unit oracles from GUI test cases," in *Proceedings of the IEEE International Conference on Software Testing, Verification and Validation (ICST)*, 2016, pp. 330–340.
- [19] M. Dhok, M. K. Ramanathan, and N. Sinha, "Type-aware concolic testing of JavaScript programs," in *Proceedings of the 38th International Conference on Software Engineering (ICSE)*, 2016, pp. 168–179.
- [20] A. Mesbah, "Software analysis for the web: Achievements and prospects," in *Proceedings of the IEEE International Conference on Software Analysis, Evolution, and Reengineering (SANER) – FoSE Track (invited)*. IEEE, 2016, pp. 91–103.
- [21] E. Andreasen, L. Gong, A. Møller, M. Pradel, M. Selakovic, K. Sen, and C. Staicu, "A survey of dynamic analysis and test generation for JavaScript," *ACM Comput. Surv.*, vol. 50, no. 5, pp. 66:1–66:36, 2017.
- [22] C. Sung, M. Kusano, N. Sinha, and C. Wang, "Static DOM event dependency analysis for testing web applications," in *Proceedings of the 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering (FSE)*, 2016, pp. 447–459.
- [23] A. Arcuri, "Longer is better: On the role of test sequence length in software testing," in *Proofceedings of the 3rd International Conference on Software Testing, Verification and Validation (ICST)*, 2010, pp. 469–478.
- [24] J. H. Andrews, A. Groce, M. Weston, and R. Xu, "Random test run length and effectiveness," in *Proofceedings of th 23rd IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 2008, pp. 19–28.
- [25] G. Fraser and A. Gargantini, "Experiments on the test case length in specification based test case generation," in *Proceedings of the 4th International Workshop on Automation of Software Test (AST)*, 2009, pp. 18–26.
- [26] Q. Xie and A. M. Memon, "Studying the characteristics of a "good" GUI test suite," in *Proceedings of the 17th International Symposium on Software Reliability Engineering (ISSRE)*, 2006, pp. 159–168.
- [27] S. Carino and J. H. Andrews, "Evaluating the effect of test case length on GUI test suite performance," in *Proceedings of the 10th IEEE/ACM International Workshop on Automation of Software Test (AST)*, 2015, pp. 13–17.
- [28] Y. Li, P. K. Das, and D. L. Dowe, "Two decades of web application testing - A survey of recent advances," *Inf. Syst.*, vol. 43, pp. 20–54, 2014.
- [29] T. Su, G. Meng, Y. Chen, K. Wu, W. Yang, Y. Yao, G. Pu, Y. Liu, and Z. Su, "Guided, stochastic model-based gui testing of android apps," in *Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering (ESEC/FSE)*, 2017, pp. 245–256.
- [30] L. Cheng, Z. Yang, and C. Wang, "Systematic reduction of GUI test sequences," in *Proceedings of the 32nd IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 2017, pp. 849–860.
- [31] P. Godefroid, *Partial-Order Methods for the Verification of Concurrent Systems - An Approach to the State-Explosion Problem*, ser. Lecture Notes in Computer Science. Springer, 1996, vol. 1032.
- [32] A. Arcuri and L. C. Briand, "A hitchhiker's guide to statistical tests for assessing randomized algorithms in software engineering," *Softw. Test., Verif. Reliab.*, vol. 24, no. 3, pp. 219–250, 2014.
- [33] A. C. D. Neto and G. H. Travassos, "A picture from the model-based testing area: Concepts, techniques, and challenges," *Advances in Computers*, vol. 80, pp. 45–120, 2010.

- [34] M. Utting, A. Pretschner, and B. Legeard, "A taxonomy of model-based testing approaches," *Softw. Test., Verif. Reliab.*, vol. 22, no. 5, pp. 297–312, 2012.
- [35] W. Li, F. Le Gall, and N. Spaseski, "A survey on model-based testing tools for test case generation," in *Proceedings of the 4th International Conference on Tools and Methods of Program Analysis (TMPA)*, 2018, pp. 77–89.
- [36] C. S. Jensen, M. R. Prasad, and A. Møller, "Automated testing with targeted event sequence generation," in *Proceedings of the International Symposium on Software Testing and Analysis (ISSTA)*, 2013, pp. 67–77.
- [37] T. Takala, M. Katara, and J. Harty, "Experiences of system-level model-based GUI testing of an android application," in *Proceedings of the 4th IEEE International Conference on Software Testing, Verification and Validation (ICST)*, 2011, pp. 377–386.
- [38] D. Amalfitano, A. R. Fasolino, P. Tramontana, B. D. Ta, and A. M. Memon, "Mobiguitar: Automated model-based testing of mobile apps," *IEEE Software*, vol. 32, no. 5, pp. 53–59, 2015.
- [39] W. Yang, M. R. Prasad, and T. Xie, "A grey-box approach for automated gui-model generation of mobile applications," in *Proceedings of the International Conference on Fundamental Approaches to Software Engineering*, 2013, pp. 250–265.
- [40] Y. M. Baek and D. H. Bae, "Automated model-based Android GUI testing using multi-level gui comparison criteria," in *Proceedings of the IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 2016, pp. 238–249.
- [41] A. Arcuri, "A theoretical and empirical analysis of the role of test sequence length in software testing for structural coverage," *IEEE Trans. Software Eng.*, vol. 38, no. 3, pp. 497–519, 2012.
- [42] A. Rau, J. Hotzkow, and A. Zeller, "Efficient GUI test generation by learning from tests of other apps," in *Proceedings of the 40th International Conference on Software Engineering: Companion Proceedings (ICSE)*, 2018, pp. 370–371.
- [43] F. Ricca and P. Tonella, "Analysis and testing of web applications," in *Proceedings of the 23rd International Conference on Software Engineering (ICSE)*, 2001, pp. 25–34.
- [44] A. A. Andrews, J. Offutt, and R. T. Alexander, "Testing web applications by modeling with fsms," *Software and System Modeling*, vol. 4, no. 3, pp. 326–345, 2005.
- [45] A. Mesbah, E. Bozdog, and A. van Deursen, "Crawling AJAX by inferring user interface state changes," in *Proceedings of the 8th International Conference on Web Engineering (ICWE)*, 2008, pp. 122–134.
- [46] A. Marchetto, P. Tonella, and F. Ricca, "State-based testing of ajax web applications," in *Proceedings of the International Symposium on Search Based Software Engineering*, 2009, pp. 3–12.
- [47] V. Dallmeier, M. Burger, T. Orth, and A. Zeller, "Webmate: a tool for testing web 2.0 applications," in *Proceedings of the Workshop on JavaScript Tools*. ACM, 2012, pp. 11–15.
- [48] N. Alshahwan and M. Harman, "Automated web application testing using search based software engineering," in *Proceedings of the 26th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 2011, pp. 3–12.
- [49] G. Wassermann, D. Yu, A. Chander, D. Dhurjati, H. Inamura, and Z. Su, "Dynamic test input generation for web applications," in *Proceedings of the ACM/SIGSOFT International Symposium on Software Testing and Analysis (ISSTA)*, 2008, pp. 249–260.
- [50] S. Artzi, A. Kiezun, J. Dolby, F. Tip, D. Dig, A. M. Paradkar, and M. D. Ernst, "Finding bugs in dynamic web applications," in *Proceedings of the ACM/SIGSOFT International Symposium on Software Testing and Analysis (ISSTA)*, 2008, pp. 261–272.
- [51] F. Gross, G. Fraser, and A. Zeller, "EXSYST: search-based GUI testing," in *Proceedings of the 34th International Conference on Software Engineering (ICSE)*, 2012, pp. 1423–1426.
- [52] J. Thomé, A. Gorla, and A. Zeller, "Search-based security testing of web applications," in *Proceedings of the 7th International Workshop on Search-Based Software Testing (SBST)*, 2014, pp. 5–14.
- [53] A. Zeller, "Search-based testing and system testing: A marriage in heaven," in *Proceedings of the 10th IEEE/ACM International Workshop on Search-Based Software Testing, (SBST@ICSE)*, 2017, pp. 49–50.
- [54] C. Nguyen, H. Yoshida, M. R. Prasad, I. Ghosh, and K. Sen, "Generating succinct test cases using don't care analysis," in *Proceedings of the 8th IEEE International Conference on Software Testing, Verification and Validation (ICST)*, 2015, pp. 1–10.
- [55] A. Rau, J. Hotzkow, and A. Zeller, "Transferring tests across web applications," in *Proceedings of the 18th International Conference on Web Engineering (ICWE)*, 2018, pp. 50–64.