

# Unit Test Generation for JavaScript

Shabnam Mirshokraie, Ali Mesbah, and Karthik Pattabiraman

*Department of Electrical and Computer Engineering, University of British Columbia, Vancouver, BC, Canada*

## SUMMARY

The highly dynamic nature of JavaScript makes it challenging to test JavaScript-based applications. Current web test automation techniques target the generation of event sequences, but they ignore testing the JavaScript code at unit level. Further they either ignore the oracle problem or simplify it through generic soft oracles. We present a framework to automatically generate test cases for JavaScript applications at two complementary levels, namely events and individual JavaScript functions. Our approach employs a combination of function coverage maximization and function state reduction algorithms to efficiently generate test cases. These test cases are strengthened by automatically generated mutation-based oracles. We empirically evaluate our approach, called JSEFT, to assess its efficacy. The results, on 13 JavaScript-based applications, show that JSEFT can detect injected faults with a high accuracy (100% precision, 72% recall). We also find that JSEFT outperforms an existing JavaScript test automation framework in terms of coverage and detected faults.

Copyright © 0000 John Wiley & Sons, Ltd.

Received ...

KEY WORDS: Test generation; oracles; JavaScript; DOM

## 1. INTRODUCTION

JavaScript plays a prominent role in modern web applications. To test their JavaScript applications, developers often write test cases using web testing frameworks such as SELENIUM (GUI tests) and QUNIT (JavaScript unit tests). Although such frameworks help to automate test execution, the test cases still need to be written manually, which is tedious and time-consuming.

Further, the event-driven nature of JavaScript as well as its runtime interaction with the Document Object Model (DOM) make JavaScript applications error-prone [27] and difficult to test. The highly dynamic behaviour of JavaScript can result in prohibitively large test suites containing thousands of test cases. This hinders the test suite comprehension. Moreover, as the size of the test suite increases, it can become expensive to execute on new versions of the application as large test suites often contain fragile, or redundant test cases [16]. However, removing test cases may negatively affect the fault finding capability of the test suite.

Researchers have recently developed automated test generation techniques for JavaScript-based applications [7, 20, 21, 23, 31]. However, current web test generation techniques suffer from two main shortcomings, namely, they:

1. Target the generation of *event sequences*, which operate at the event-level or DOM-level to cover the state space of the application. These techniques fail to capture faults that do not propagate to an observable DOM state. As such, they potentially miss this portion of code-level JavaScript faults. In order to capture such faults, effective test generation techniques need to target the code at the JavaScript unit-level, in addition to the event-level.
2. Either ignore the oracle problem altogether or simplify it through generic *soft oracles*, such as W3C HTML validation [7, 23], or JavaScript runtime exceptions [7]. A generated test

case without assertions is not useful since coverage alone is not the goal of software testing. For such generated test cases, the tester still needs to manually write many assertions, which is time and effort intensive. On the other hand, soft oracles target generic fault types and are limited in their fault finding capabilities. However, to be practically useful, unit testing requires strong oracles to determine whether the application under test executes correctly.

To address these two shortcomings, we propose an automated test case generation technique for JavaScript applications.

Our approach, called JSEFT (JavaScript Event and Function Testing) operates through a three step process. First, it dynamically explores the event-space of the application using a *function coverage maximization* method, to infer a test model. Then, it generates test cases at two complementary levels, namely, DOM event and JavaScript functions. Our technique employs a novel *function state reduction* algorithm to minimize the number of function-level states needed for test generation. Finally, it automatically generates test oracles, through a mutation-based algorithm. This paper is a substantially expanded and revised version of our paper from early 2015 [26] (see appendix for differences from that paper).

This work makes the following main contributions:

- An automatic technique to generate test cases for JavaScript functions and events.
- A combination of *function coverage maximization* and *function state reduction* algorithms to efficiently generate unit test cases;
- A mutation-based algorithm to effectively generate test oracles, capable of detecting regression JavaScript and DOM-level faults;
- The implementation of our technique in a tool called JSEFT, which is publicly available [3];
- An empirical evaluation to assess the efficacy of JSEFT using 13 JavaScript applications.

The results of our evaluation show that on average (1) the generated test suite by JSEFT achieves a 68% JavaScript code coverage, (2) compared to ARTEMIS, a feedback-directed JavaScript testing framework [7], JSEFT achieves 53% better coverage, and (3) the test oracles generated by JSEFT are able to detect injected faults with 100% precision and 72% recall.

## 2. CHALLENGES AND MOTIVATION

In this section, we illustrate some of the challenges associated with test generation for JavaScript applications.

Figure 1 presents a snippet of a JavaScript game application that we use as a running example throughout the paper. This simple example uses the popular jQuery library [1] and contains four main JavaScript functions:

1. `cellClicked` is bound to the event-handlers of DOM elements with IDs `cell0` and `cell1` (Lines 34–35). These two DOM elements become available when the DOM is fully loaded (Line 32). Depending on the element clicked, `cellClicked` inserts a `div` element with ID `divElem` (Line 3) after the clicked element and makes it clickable by attaching either `setup` or `setDim` as its event-handler function (Lines 5–6, 9–10).
2. `setup` calls `setDim` (Line 15) to change the value of the global variable `currentDim`. It further makes an element with ID `startCell` clickable by setting its event-handler to `start` (Line 16).
3. `setDim` receives an input variable. It performs some computations to set the `height` value of the `css` property of a DOM element with ID `endCell` and the value of `currentDim` (Lines 20–22). It also returns the computed dimension.
4. `start` is called at runtime when the element with ID `startCell` is clicked (Line 16), which either updates the width dimension of the element on which it was called, or removes the element (Lines 27–29).

```

1 var currentDim=20;
2 function cellClicked() {
3   var divTag = '<div id='divElem' />';
4   if($(this).attr('id') == 'cell0'){
5     $('#cell0').after(divTag);
6     $('#div #divElem').click(setup);
7   }
8   else if($(this).attr('id') == 'cell1'){
9     $('#cell1').after(divTag);
10    $('#div #divElem').click(function(){setDim(20)});
11  }
12 }

14 function setup() {
15   setDim(10);
16   $('#startCell').click(start);
17 }

19 function setDim(dimension) {
20   var dim=($('#endCell').width() + $('#endCell').height())/dimension;
21   currentDim += dim;
22   $('#endCell').css('height', dim+'px');
23   return dim;
24 }

26 function start() {
27   if(currentDim > 40)
28     $(this).css('height', currentDim+'px');
29   else $(this).remove();
30 }

32 $document.ready(function() {
33   ...
34   $('#cell0').click(cellClicked);
35   $('#cell1').click(cellClicked);
36 });

```

Figure 1. JavaScript code of the running example.

There are four main challenges in testing JavaScript applications.

The first challenge is that a fault may not immediately propagate into a DOM-level observable failure. For example, if the '+' sign in Line 21 is mistakenly replaced by '-', the affected result does not immediately propagate to the observable DOM state after the function exits. While this mistakenly changes the value of a global variable, `currentDim`, which is later used in `start` (Line 27), it neither affects the returned value of the `setDim` function nor the `css` value of element `endCell`. Therefore, a GUI-level event-based testing approach may not help to detect the fault in this case.

The second challenge is related to fault localization; even if the fault propagates to a future DOM state and a DOM-level test case detects it, finding the actual location of the fault is challenging for the tester as the DOM-level test case is agnostic of the JavaScript code. However, a unit test case that targets individual functions, e.g., `setDim` in this running example, helps a tester to spot the fault, and thus easily resolve it.

The third challenge pertains to the event-driven dynamic nature of JavaScript, and its extensive interaction with the DOM resulting in many state permutations and execution paths. In the initial state of the example, clicking on `cell0` or `cell1` takes the browser to two different states as a result of the `if-else` statement in Lines 4 and 8 of the function `cellClicked`. Even in this simple example, expanding either of the resulting states has different consequences due to different functions that can be potentially triggered. Executing either `setup` or `setDim` in Lines 6 and 10 results in different execution paths, DOM states, and code coverage. It is this dynamic interaction of the JavaScript code with the DOM (and indirectly CSS) at runtime that makes it challenging to generate test cases for JavaScript applications.

The fourth important challenge in unit testing JavaScript functions that have DOM interactions, such as `setDim`, is that the DOM tree in the state expected by the function, has to be present during unit test execution. Otherwise the test will fail due to a `null` or `undefined` exception. This situation arises often in modern web applications that have many DOM interactions.

### 3. APPROACH

Our main goal in this work is to generate client-side test cases coupled with effective test oracles, capable of detecting regression JavaScript and DOM-level faults. Further, we aim to achieve this goal as efficiently as possible. Hence, we make two design decisions. First, we assume that there is a finite amount of time available to generate test cases. Consequently we guide the test generation to maximize coverage under a given time constraint. The second decision is to minimize the number of test cases and oracles generated to only include those that are essential in detecting potential faults. Consequently, to examine the correctness of the test suite generated, the tester would only need to examine a small set of assertions, which minimizes their effort.

Our approach generates test cases and oracles at two complementary levels:

**DOM-level event-based tests** consist of DOM-level event sequences and assertions to check the application's behaviour from an end-user's perspective.

**Function-level unit tests** consist of unit tests with assertions that verify the functionality of JavaScript code at the function level.

An overview of the technique is depicted in Figure 2. At a high level, our approach is composed of three main steps:

1. In the first step (Section 3.1), we dynamically explore various states of a given web application, in such a way as to maximize the number of functions that are covered throughout the program execution. The output of this initial step is a state-flow graph (SFG) [23], capturing the explored dynamic DOM states and event-based transitions between them.
2. In the second step (Section 3.2), we use the inferred SFG to generate event-based test cases. We run the generated tests against an instrumented version of the application. From the execution trace obtained, we extract DOM element states as well as JavaScript function states at the entry and exit points.
3. To create effective test oracles for DOM-level test cases, we perform a DOM-based mutation analysis (Section 6). The test oracles are generated at the DOM level by comparing DOM properties inferred from the original and the mutated versions of the application. To reduce the number of function-level unit test cases as well as generated oracles to only those that are constructive, we automatically generate mutated versions of the application (Section 6). We devise a *state reduction* algorithm that minimizes the number of states by selecting representative mutated function states. The corresponding original states are used to generate function-level unit tests. Assuming that the original version of the application is fault-free, the test oracles are then generated at the code level by comparing the states traced from the original and the mutated versions.

#### 3.1. Maximizing Function Coverage

In this step, our goal is to maximize the number of functions that can be covered, while exercising the program's event space. To that end, our approach combines static and dynamic analysis to decide which state and event(s) should be selected for expansion to maximize the probability of covering uncovered JavaScript functions. While exploring the web application under test, our function coverage maximization algorithm selects a next state for exploration, which has the maximum value of the sum of the following two metrics:

**1. Potential Uncovered Functions.** This pertains to the total number of unexecuted functions that can potentially be visited through the execution of DOM events in a given DOM state  $s_i$ . When

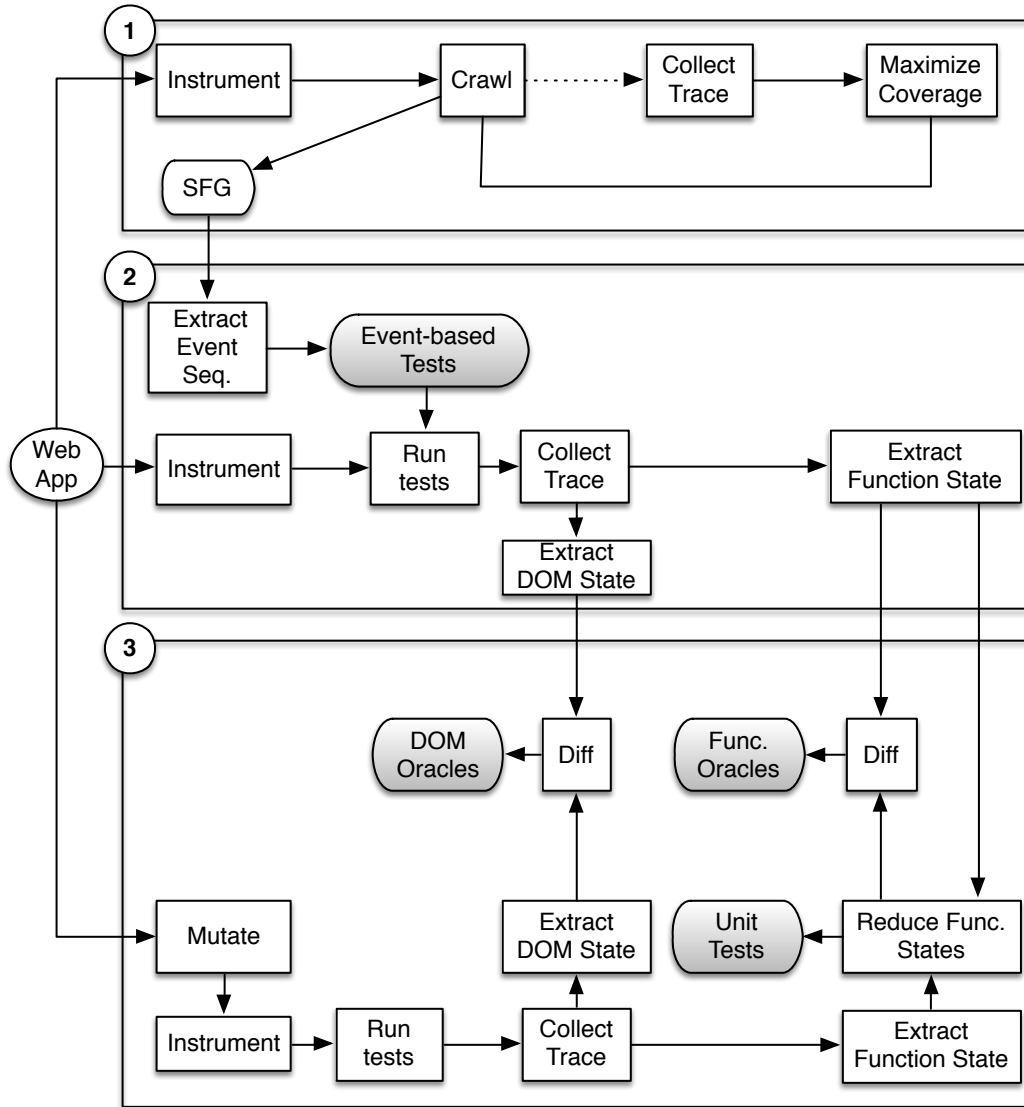


Figure 2. Overview of our test generation approach.

a given function  $f_i$  is set as the event-handler of a DOM element  $d \in s_i$ , it makes the element a potential clickable element in  $s_i$ . This can be achieved through various patterns in web applications depending on which DOM event model level is adopted. To calculate this metric, our algorithm identifies all JavaScript functions that are directly or indirectly attached to DOM elements as event handlers, in  $s_i$  through code instrumentation and execution trace monitoring.

**2. Potential Clickable Elements.** The second metric, used to select a state for expansion, pertains to the number of DOM elements that can potentially become clickable elements. If the event-handlers bound to those clickables are triggered, new (uncovered) functions will be executed. To obtain this number, we statically analyze the previously obtained *potential uncovered functions* within a given state in search of such elements.

While exploring the application, the next state for expansion is selected by adding the two metrics and choosing the state with the highest sum. The procedure repeats the aforementioned steps until the designated time limit, or state space size is reached.

In the running example of Figure 1, in the initial state, clicking on elements with IDs `cell10` and `cell11` results in two different states due to an `if-else` statement in Lines 4 and 8 of

cellClicked. Let's call the state in which a DIV element is located after the element with ID cell0 as  $s_0$ , and the state in which a DIV element is placed after the element with ID cell1 as  $s_1$ . If state  $s_0$ , with the clickable cell0, is chosen for expansion, function setup is called. As shown in Line 15, setup calls setDim, and thus, by expanding  $s_0$  both of the aforementioned functions get called by a single click. Moreover, a potential clickable element is also created in Line 16, with start as the event-handler. Therefore, expanding  $s_1$  results only in the execution of setDim, while expanding  $s_0$  results in the execution of functions setup, setDim, and a potential execution of start in future states. At the end of this step, we obtain a state-flow graph of the application that can be used in the next test generation step.

### 3.2. Extracting DOM and JavaScript Function States

In the second step, our technique extracts sequences of events from the inferred state-flow graph. These sequences of events are used in our test case generation process.

**DOM-level event-based testing.** To verify the behaviour of the application at the user interface level, each event path, taken from the initial state (Index) to a leaf node in the state-flow graph, is used to generate DOM event-based test cases. Each extracted path is converted into a JUNIT SELENIUM-based test case, which executes the sequence of events, starting from the initial DOM state. Going back to our running example, one possible event sequence to generate is: `$('#cell0').click→$('div #divElem').click→$('#startCell').click`.

To collect the required trace data, we capture all DOM elements and their attributes after each event in the test path is fired. This trace is later used in our DOM oracle comparison, as explained in Section 6.

**Extracting JavaScript function states.** To generate unit tests that target JavaScript functions directly (as opposed to event-triggered function executions), we log the state of each function at their entry and exit point, during execution. To that end, we instrument the code to trace various entities. At the entry point of a given JavaScript function we collect (1) function parameters including passed variables, objects, functions, and DOM elements, (2) global variables used in the function, and (3) the current DOM structure just before the function is executed. At the exit point of the JavaScript function and before every return statement, we log the state of the (1) return value of the function, (2) global variables that have been accessed in that function, and (3) DOM elements accessed (read/written) in the function. At each of the above points, our instrumentation records the name, runtime type, and actual values. The dynamic type is stored because JavaScript is a dynamically typed language, meaning that the variable types cannot be determined statically. Note that complex JavaScript objects can contain circular or multiple references (e.g., in JSON format). To handle such cases, we perform a de-serialization process in which we replace such references by an object in the form of  $\$ref : Path$ , where *Path* denotes a *JSONPath* string\* that indicates the target path of the reference.

In addition to function entry and exit points, we log information required for calling the function from the generated test cases. JavaScript functions that are accessible in the public scope are mainly defined in (1) the global scope directly (e.g., `function f() { ... }`), (2) variable assignments in the global scope (e.g., `var f = function() { ... }`), (3) constructor functions (e.g., `function constructor() { this.member = function() { ... } }`), and (4) prototypes (e.g., `Constructor.prototype.f = function() { ... }`). Functions in the first and second case are easy to call from test cases. For the third case, the constructor function is called via the new operator to create an object type, which can be used to access the object's properties (e.g., `container = new Constructor(); container.member();`). This allows us to access the inner function, which is a member of the constructor function in the above example. For the prototype case, the function can be invoked through `container.f()` from a test case.

\* <http://goessner.net/articles/JsonPath/>

Going back to our running example in Figure 1, at the entry point of `setDim`, we log the value and type of both the input parameter `dimension` and global variable `currentDim`, which is accessed in the function. Similarly, at the exit point, we log the values and types of the returned variable `dim` and `currentDim`.

In addition to the values logged above, we need to capture the DOM state for functions that interact with the DOM. This is to address the fourth challenge outlined in Section 2. To mitigate this problem, we capture the state of the DOM just before the function starts its execution, and include that as a *test fixture* [4] in the generated unit test case.

In the running example, at the entry point of `setDim`, we log the `innerHTML` of the current DOM as the function contains several calls to the DOM, e.g., retrieving the element with ID `endCell` in Line 22. We further include in our execution trace the way DOM elements and their attributes are modified by the JavaScript function at runtime. The information that we log for accessed DOM elements includes the ID attribute, the XPath position of the element on the DOM tree, and all the modified attributes. Collecting this information is essential for oracle generation in the next step. We use a set to keep the information about DOM modifications, so that we can record the latest changes to a DOM element without any duplication within the function. For instance, we record ID as well as both `width` and `height` properties of the `endCell` element.

Once our instrumentation is carried out, we run the generated event sequences obtained from the state-flow graph. This way, we produce an execution trace that contains:

- Information required for preparing the environment for each function to be executed in a test case, including its input parameters, used global variables, and the DOM tree in a state that is expected by the function;
- Necessary entities that need to be assessed after the function is executed, including the function's output as well as the touched DOM elements and their attributes (The actual assessment process is explained in Section 6).

### 3.3. Mutation Analysis

In the third step, we generate our function-level unit test cases. Moreover, our approach automatically generates test oracles for the DOM as well as unit level test cases, as depicted in the third step of Figure 2. Instead of randomly generating assertions, our oracle generation uses a mutation-based process.

Mutation testing is typically used to evaluate the quality of a test suite [9], or to generate test cases that kill mutants [15]. In our approach, we adopt mutation testing to (1) reduce the number of function-level unit test cases, (2) reduce the number of assertions automatically generated, and (3) target critical and error-prone portions of the application. Hence, the tester would only need to examine a small set of effective assertions to verify the correctness of the generated oracles.

In the following we describe our function state reduction technique to reduce the number of generated unit test cases, followed by our approach for generating oracles at DOM and unit level test cases.

**3.3.1. JavaScript function-level unit testing** As mentioned in Section 2, the highly dynamic nature of JavaScript applications can result in a huge number of function states. Capturing all these different states for the purpose of test case generation can impede the test suite's comprehension. Therefore, we apply a mutation-based function state reduction method to minimize the number of function states needed for test generation. Note that these mutations are later used to create our unit level test oracles.

**Function state reduction.** The reduction technique is based on classification of mutated function states according to the impact of the modification on the function's behaviour. A mutated state is the affected function state as a result of an injected mutation (details of the injected mutations are explained in Section 6). Mutated states that are either mutually equivalent or have similar characteristics can be discarded. Although they are not equivalent to the original version of the

**Algorithm 1:** Function State Reduction

---

**input** :  $\{origSt, mutSt | Mutation(f) : origSt \rightarrow mutSt, origSt_i \in OST_f, mutSt_i \in MST_f\}$ , where  $origSt$  is the set of original function states, and  $mutSt_i$  is the set of mutated function states for a given function  $f$

**output**: The obtained reduced states set  $reducedStates$

**begin**

```

1  for  $mutSt_i \in MST_f$  do
2     $L = 1; MStSet_L \leftarrow \emptyset$ 
3    if  $BRNCovLNS[mutSt_i] \neq BRNCovLNS[MStSet]_{l=1}^L$  then
4       $MStSet_{L+1} \leftarrow mutSt_i$ 
5       $L++$ 
6    else
7       $MStSet_l \leftarrow mutSt_i \cup MStSet_l$ 
8     $K = L + 1; MStSet_K \leftarrow \emptyset$ 
9    if  $DOMPROPS[mutSt_i] \neq DOMPROPS[MStSet]_{k=L+1}^K \parallel$ 
       $RETPROP[mutSt_i] \neq RETPROP[MStSet]_{k=L+1}^K$  then
10      $MStSet_{K+1} \leftarrow mutSt_i$ 
11      $K++$ 
12   else
13      $MStSet_k \leftarrow mutSt_i \cup MStSet_k$ 
14   while  $MStSet_{K+L} \neq \emptyset$  do
15      $SelectedSt \leftarrow SELECTMAXST(mutSt_i | mutSt_i \cap MStSet_{j=1}^{K+L})$ 
16      $reducMStates.ADD(SelectedSt)$ 
17      $MStSet_{K+L} \leftarrow MStSet_{K+L} - SelectedSt$ 
18   for  $origSt_i \in OST_f$  do
19      $M = 1$ 
20      $origStSet_M \leftarrow \emptyset$ 
21     if  $BRNCovLNS[origSt_i] \neq BRNCovLNS[origStSet]_{m=1}^M$  then
22        $origStSet_{M+1} \leftarrow origSt_i$ 
23        $M++$ 
24     else
25        $origStSet_m \leftarrow origSt_i \cup origStSet_m$ 
26    $reducedStates \leftarrow \{corrOrigSt | Mutation(f) : corrOrigSt \rightarrow reducMStates, corrOrigSt_i \in OST_f\}$ 
27   for  $oSt_i \in origStSet$  do
28     if  $oSt_i \cap reducedStates = \emptyset$  then
29        $reducedStates \leftarrow SelectRand(oSt_i) \cup reducedStates$ 
30   return  $reducedStates$ 

```

---

program, there would be less point in including all of them as the test case that can detect a given mutation is potentially able to detect mutations with similar impacts as well.

To categorize mutated states we assess the impact degree of a mutation on the function state in terms of covered branches within the function, the characteristics of the function's return variable/object as well as the accessed DOM elements.

**Branch coverage:** Taking different branches in a given function can change its behaviour. Thus, mutated states that result in a different covered branch should be taken into account while generating test cases. Going back to our example in Figure 1, executing either of the branches in lines 27 and 29 clearly takes the application into a different DOM state. In this example, we need to include the mutated states of the `start` function that result in different covered branches, e.g., two different mutated function states where the value of the global variable `currentDim` falls into different boundaries.

**Return variable/object characteristics:** Properties of an object can alter in JavaScript at runtime. Moreover, variable's type can dynamically change. This can result in changes in the expected outcome of the function. Going back to our example, if `dim` is mutated such that it is assigned a `string` value before adding it to `currentDim` (Line 21) in function `setDim`,



the returned value of the function becomes the `string` concatenation of the two values rather than the expected numerical addition.

**Accessed DOM properties:** Changes in DOM elements and their properties accessed in a function can affect the behaviour of the function. For example, in line 29 `this` keyword refers to the clicked DOM element of which function `start` is an event-handler. Assuming that  $\text{currentDim} \leq 40$ , depending on which DOM element is clicked, by removing the element in line 29 the resulting state of the function `start` differs. If the mutation alters the element that invokes the function `start`, the expected outcome of the function changes as well. Therefore, we take into consideration the DOM elements accessed by the function as well as the type of accessed DOM properties.

One issue with removing similar mutated states is that the code coverage can be lowered when the test cases are constructed from the corresponding original states. To make sure that our reduction technique does not lower the coverage, we measure the branch coverage achieved by each of the original function states. Function states with equal branch coverage are categorized under the same set, and the reduced set of mutated states is compared against the categorized original states. The final set includes at least one function state from each of the covered branch original sets.

Algorithm 1 shows our function state reduction algorithm. The algorithm first collects covered branches of individual functions per mutated state ( $\text{BRNCovLNS}[\text{mutSt}_i]$  in Line 3). Each mutated function's states exhibiting same covered branches are categorized under the same set of states (Lines 4 and 7).  $\text{MStSet}_l$  corresponds to the set of mutated function states, which are classified according to their covered branches, where  $l = 1, \dots, L$  and  $L$  is the number of current classified sets in covered branch category. Similarly, affected function states with the same accessed DOM characteristics as well as return variable/object properties, are put into the same set of states (Lines 10 and 13).  $\text{MStSet}_k$  corresponds to the set of mutated function states, which are classified according to their DOM/return properties, where  $k = 1, \dots, K$  and  $K$  is the number of current classified sets in that category. After classifying each mutated function's states into several sets, we cover each set by selecting one of its common states. The state selection step is a *set cover problem* [8], i.e., given a universe  $U$  and a family  $S$  of subsets of  $U$ , a cover is a subfamily  $C \subseteq S$  of sets whose union is  $U$ . Sets to be covered in our algorithm are  $\text{MStSet}_{K+L}$ , where  $\text{mutSt}_i \in \text{MStSet}_{K+L}$ . We use a common greedy algorithm for obtaining the minimum number of states that can cover all the possible sets (Lines 15-17). To this end,  $\text{reducMStates}$  (Line 16) contains reduced set of mutated states. Note that there is a corresponding original state for each of the mutated function states in the reduced set ( $\text{corrOrigSt}$  in Line 26). As mentioned earlier, to prevent from negative effect of the state reduction technique on the original code coverage, we further classify original function states according to the covered branches (Lines 18-25).  $\text{origStSet}_m$  contains the set of categorized function states, where  $m = 1, \dots, M$  and  $M$  is the number of current classified sets. The final reduced list of original function states ( $\text{reducedStates}$ ) contains at least one state from each of the covered branches category (Lines 27-29). The final list of states is returned in Line 30.

**3.3.2. Oracle Generation** Algorithm 2 shows our algorithm for generating test oracles. At a high level, the technique iteratively executes the following steps:

1. A mutant is created by injecting a single fault into the original version of the web application (Line 9 and 19 in Algorithm 2 for DOM mutation and code-level mutation, respectively),
2. Related entry/exit program states at the DOM and JavaScript function levels of the mutant and the original version are captured.  $\text{OnEvDomSt}$  in Line 4 is the original DOM state on which the event  $Ev$  is triggered,  $\text{AfterEvDomSt}$  in line 5 is the observed DOM state after the event is triggered,  $\text{MutDom}$  in line 9 is the mutated DOM, and  $\text{ChangedSt}$  in line 10 is the corresponding affected state for DOM mutations.  $\text{FcExit}$  in Line 22 is the exit state of the function in the original application and  $\text{MutFcExit}$  in line 23 is the corresponding exit state for that function after the application is mutated for function-level mutations.

**Algorithm 2:** Oracle Generation

---

**input** : A Web application (*App*), list of event sequences obtained from SFG (*EvSeq*), maximum number of mutations (*n*)

**output**: Assertions for function-level (*FcAsserts*) and DOM event-level tests (*DomAsserts*)

---

```

1  App ← INSTRUMENT(App)
   begin
2      while GenMuts < n do
3          foreach EvSeq ∈ SFG do
4              OnEvDomSt ← Trace.GETONEVDOMST(Ev ∈ EvSeq)
5              AfterEvDomSt ← Trace.GETAFTEREVDOMST(Ev ∈ EvSeq)
6              AccdDomProps ← GETACCDOMNDS(OnEvDomSt)
7              EquivalentDomMut ← true
8              while EquivalentDomMut do
9                  MutDom ← MUTATEDOM(AccdDomProps, OnEvDomSt)
10                 ChangedSt ← EvSeq.EXECEVENT(MutDom)
11                 DiffChangedSt, AfterEvDomSt ← DIFF(ChangedSt, AfterEvDomSt)
12                 if DiffChangedSt, AfterEvDomSt ≠ ∅ then
13                     EquivalentDomMut ← false
14                     DomAsserti = DiffChangedSt, AfterEvDomSt
15                     DomAssertsEv, AfterEvDomSt = ∩ DomAsserti
16
17                 RedFcSts ← Trace.GETREDFCSTS()
18                 EquivalentCodeMut ← true
19                 while EquivalentCodeMut do
20                     MutApp ← MUTATEJSCODE(App)
21                     MutFcSts ← EvSeq.EXECEVENT(MutApp)
22                     foreach FcEntry ∈ RedFcSts.GETFCENTRIES do
23                         FcExit ← RedFcSts.GETFCEXIT(FcEntry)
24                         MutFcExit ← MutFcSts.GETMUTFCEXIT(FcEntry)
25                         DiffFcExit, MutFcExit ← DIFF(FcExit, MutFcExit)
26                         if DiffFcExit, MutFcExit ≠ ∅ then
27                             EquivalentCodeMut ← false
28                             FcAsserti = ∩ DiffFcExit, MutFcExit
29                             FcAssertsFcEntry = ∪ FcAsserti
29
30     return {FcAsserts, DOMAsserts}

```

---

3. Relevant observed state differences at each level are detected and abstracted into test oracles (DIFF in Line 11 and 24 for DOM and function-level oracles, respectively),
4. The generated assertions (Lines 15 and 28) are injected into the corresponding test cases.

**DOM-level event-based test oracles.** After an event is triggered in the generated SELENIUM test case, the resulting DOM state needs to be compared against the expected structure. One naive approach would be to compare the DOM tree in its entirety, after the event execution. Not only is this approach inefficient, it results in brittle test-cases, i.e., the smallest update on the user interface can break the test case. We propose an alternative approach that utilizes *DOM mutation testing* to detect and selectively compare only those DOM elements and attributes that are affected by an injected fault at the DOM-level of the application. Our DOM mutations target only the elements that have been accessed (read/written) during execution, and thus have a larger impact on the application's behaviour. To select proper DOM elements for mutation, we instrument JavaScript functions that interact with the DOM, i.e., code that either accesses or modifies DOM elements.

We execute the instrumented application by running the generated SELENIUM test cases and record each accessed DOM element, its attributes, the triggered event on the DOM state, and the DOM state after the event is triggered (GETONEVDOMST in line 4, GETAFTEREVDOMST in line 5, and GETACCDOMNDS in line 6 to retrieve the original DOM state, DOM state after event *Ev* is triggered, and the accessed DOM properties as event *Ev* is triggered, respectively, in Algorithm 2). To perform the actual mutation, as the application is re-executed using the same sequence of events, we mutate the recorded DOM elements, one at a time, before the corresponding

```

1 @Test
2 public void testCase1() {
3     WebElement divElem=driver.findElements(By.id("divElem"));
4     divElem.click();
5     int endCellHeight=driver.findElements(By.id("endCell")).getSize().height;
6     assertEquals(endCellHeight, 30);
7     WebElement startCell=driver.findElements(By.id("startCell"));
8     startCell.click();
9     boolean exists=driver.findElements(By.id("startCell")).size!=0;
10    assertTrue(exists);
11    int startCellHeight=driver.findElements(By.id("startCell")).getSize().height;
12    assertEquals(startCellHeight, 50);
13 }

```

Figure 3. Generated SELENIUM test case.

event is fired. *MUTATEDOM* in line 9 mutates the DOM elements, and *EvSeq.EXECEVENT* in line 10 executes the event sequence on the mutated DOM. The mutation operators include (1) deleting a DOM element, and (2) changing the attribute, accessed during the original execution. As we mutate the DOM, we collect the current state of DOM elements and attributes.

Figure 3 shows part of a DOM-level test case generated for the running example. Going back to our running example, as a result of clicking on  $\$(\text{'div \#divElem'})$  in our previously obtained event sequence  $\$(\text{'\#cell0'}) \rightarrow \text{click} \rightarrow \$(\text{'div \#divElem'}) \rightarrow \text{click} \rightarrow \$(\text{'\#startCell'})$ , the height and width properties of DOM element with ID *endCell*, and the DOM element with ID *startCell* are accessed. One possible DOM mutation is altering the width value of the *endCell* element before click on  $\$(\text{'div \#divElem'})$  happens. We log the consequences of this modification after the click event on  $\$(\text{'div \#divElem'})$  as well as the remaining events. This mutation affects the height property of DOM element with ID *endCell* in the resulting DOM state from clicking on  $\$(\text{'div \#divElem'})$ . Line 6 in Figure 3 shows the corresponding assertion. Furthermore, Assuming that the DOM mutation makes  $\text{currentDim} \leq 40$  in line 27, after click on element *\#startCell* happens, the element is removed and no longer exists in the resulting DOM state. The generated assertion is shown in line 10 of Figure 3.

Hence, we obtain two sets of execution traces that contain information about the state of DOM elements for each fired event in the original and mutated application. By comparing these two traces (*DIFF* in line 11 in Algorithm 2), we identify all changed DOM elements and generate assertions for these elements. Note that any changes detected by the *DIFF* operator (line 12 in Algorithm 2) is an indication that the corresponding DOM mutation is *not equivalent* (line 13); if no change is detected, another DOM mutation is generated. We automatically place the generated assertion immediately after the corresponding line of code that executed the event, in the generated event-based (SELENIUM) test case. *DomAssertsEv,AfterEvDomSt* in line 15 contains all DOM assertions for the state *AfterEvDOMSt* and the triggered event *Ev*.

**Function-level test oracles.** To seed code level faults, we use our recently developed JavaScript mutation testing tool, *MUTANDIS* [25]. Mutations generated by *MUTANDIS* are selected through a *function rank* metric, which ranks functions in terms of their relative importance from the application's behaviour point of view. The mutation operators are chosen from a list of common operators, such as changing the value of a variable or modifying a conditional statement. Once a mutant is produced (*MUTATEJSCODE* in line 19), it is automatically instrumented. We collect a new execution trace from the mutated program by executing the same sequence of events that was used on the original version of the application. This way, the state of each JavaScript function is extracted at its entry and exit points. *RedFcSts.GETFCENTRIES* in line 21 retrieves the function's entries from the reduced function's states. *GETFCEXIT* in line 22, and *GETMUTFCEXIT* in line 23 retrieve the corresponding function's exit state in the original and mutated application. This process is similar to the function state extraction algorithm explained in Section 3.2.

```

1 test("Testing setDim",4,function(){
2   var fixture = $("#qunit-fixture");
3   fixture.append("<button id=\"cell0\"> <div id=\"divElem\"/> </button> <div id=\"endCell\" style=\"height:200px;width:100px;\"/>");
4   var currentDim=20;
5   var result= setDim(10);
6   equal(result, 30);
7   equal(currentDim, 50);
8   ok($("#endCell").length > 0);
9   equal($("#endCell").css('height'), 30); });

```

Figure 4. Generated QUNIT test case.

After the execution traces are collected for all the generated mutants, we generate function-level test oracles by comparing the execution trace of the original application with the traces we obtained from the modified versions (DIFF in line 24 in Algorithm 2). If the DIFF operator detects no changes (line 25 of the algorithm), an equivalent mutant is detected, and thus another mutant will be generated.

Our function-level oracle generation targets *postcondition assertions*. Such postcondition assertions can be used to examine the expected behaviour of a given function after it is executed in a unit test case. Our technique generates postcondition assertions for all functions that exhibit a *different exit-point state* but the *same entry-point state*, in the mutated execution traces.  $FcAssert_i$  in line 27 contains all such post condition assertions. Due to the dynamic and asynchronous behaviour of JavaScript applications, a function with the same entry state can exhibit different outputs when called multiple times. In this case, we need to combine assertions to make sure that the generated test cases do not mistakenly fail.  $FcAsserts_{FcEntry}$  in line 28 contains the union of function assertions generated for the same entry but different outputs during multiple executions. Let's consider a function  $f$  with an entry state  $entry$  in the original version of the application ( $A$ ), with two different exit states  $exit_1$  and  $exit_2$ . If in the mutated version of the application ( $A_m$ ),  $f$  exhibits an exit state  $exit_m$  that is different from both  $exit_1$  and  $exit_2$ , then we combine the resulting assertions as follows:  $assert1(exit_1, expRes_1) || assert2(exit_2, expRes_2)$ , where the expected values  $expRes_1$  and  $expRes_2$  are obtained from the execution trace of  $A$ .

Each assertion for a function contains (1) the function's returned value, (2) the used global variables in that function, and/or (3) the accessed DOM element in that function. Each assertion is coupled with the expected value obtained from the execution trace of the original version.

The generated assertions that target variables, compare the value as well as the runtime type against the expected ones. An oracle that targets a DOM element, first checks the existence of that DOM element. If the element exists, it checks the attributes of the element by comparing them against the observed values in the original execution trace. Assuming that `width` and `height` are 100 and 200 accordingly in Figure 1, and '+' sign is mutated to '-' in line 20 of the running example in Figure 1, the mutation affects the global variable `currentDim`, `height` property of element with ID `endCell`, and the returned value of the function `setDim`. Figure 4 shows a QUNIT test case for `setDim` function according to this mutation with the generated assertions.

### 3.4. Tool Implementation

We have implemented our JavaScript test and oracle generation approach in an automated tool called JSEFT. The tool is written in Java and is publicly available for download [3]. Our implementation requires no browser modifications, and is hence portable. For JavaScript code interception, we use a web proxy, which enables us to automatically instrument JavaScript code before it reaches the browser. The crawler for JSEFT extends and builds on top of the event-based crawler, CRAWLJAX [22], with random input generation enabled for form inputs. As mentioned before, to mutate JavaScript code, we use our recently developed mutation testing tool, MUTANDIS [25]. The upper-bound for the number of mutations can be specified by the user. However, the default is 50 for code-level and 20 for DOM-level mutations. We observed that these default numbers provide a balanced trade-off between oracle generation time, and the fault finding capability of the tool. DOM-level test

Table I. Characteristics of the experimental objects.

ID	Name	LOC	URL
1	SameGame	206	<a href="http://crawljax.com/same-game/">crawljax.com/same-game/</a>
2	Tunnel	334	<a href="http://arcade.christianmontoya.com/tunnel/">arcade.christianmontoya.com/tunnel/</a>
3	GhostBusters	282	<a href="http://10k.aneventapart.com/2/Uploads/657/">10k.aneventapart.com/2/Uploads/657/</a>
4	Peg	509	<a href="http://www.cccontheweb.org/peggame.htm">www.cccontheweb.org/peggame.htm</a>
5	BunnyHunt	580	<a href="http://www.themaninblue.com/experiment/BunnyHunt/">http://www.themaninblue.com/experiment/BunnyHunt/</a>
6	AjaxTabs	592	<a href="https://github.com/amazingSurge/jquery-tabs/">https://github.com/amazingSurge/jquery-tabs/</a>
7	NarrowDesign	1,005	<a href="http://www.narrowdesign.com">http://www.narrowdesign.com</a>
8	JointLondon	1,211	<a href="http://www.jointlondon.com">http://www.jointlondon.com</a>
9	FractalViewer	1,245	<a href="http://onecm.com/projects/canopy/">onecm.com/projects/canopy/</a>
10	SimpleCart	1,900	<a href="https://github.com/wojodesign/simplecart-js/">https://github.com/wojodesign/simplecart-js/</a>
11	WymEditor	3,035	<a href="http://www.wymeditor.org">http://www.wymeditor.org</a>
12	TuduList	1,963	<a href="http://tudu.ess.ch/tudu">http://tudu.ess.ch/tudu</a>
13	TinyMCE	26,908	<a href="http://www.tinymce.com">http://www.tinymce.com</a>

cases are generated in a JUNIT format that uses SELENIUM (WebDriver) APIs to fire events on the application's DOM inside the browser. JavaScript function-level tests are generated in the QUNIT unit testing framework [4], capable of testing any generic JavaScript code.

#### 4. EMPIRICAL EVALUATION

To quantitatively assess the efficacy of our test generation approach, we have conducted an empirical study, in which we address the following research questions:

**RQ1** How effective is JSEFT in generating test cases with high coverage?

**RQ2** How capable is JSEFT of generating test oracles that detect regression faults?

**RQ3** How effective is the state reduction technique in reducing the number of function states?

**RQ4** How does JSEFT compare to existing automated JavaScript testing frameworks?

JSEFT and all our experimental data in this paper are available for download [3].

##### 4.1. Objects

Our study includes thirteen JavaScript-based applications in total. Table I presents each application's ID, name, lines of custom JavaScript code (LOC, excluding JavaScript libraries) and resource. The first five are web-based games. AjaxTabs is a JQUERY plugin for creating tabs. NarrowDesign and JointLondon are websites. FractalViewer is a fractal tree zoom application. SimpleCart is a shopping cart library, WymEditor is a web-based HTML editor, TuduList is a web-based task management application, and TinyMCE is a JavaScript based WYSIWYG editor control. The applications range from 206 to 27K lines of JavaScript code.

The experimental objects are open-source and cover different application types. All the applications are interactive in nature and extensively use JavaScript on the client-side.

##### 4.2. Setup

To address our research questions, we provide the URL of each experimental object to JSEFT. Test cases are then automatically generated by JSEFT. We give JSEFT 10 minutes in total for each application. 5 minutes of the total time is designated for the dynamic exploration step.

**Test Case Generation (RQ1).** To measure client-side code coverage, we use JSCOVER [2], an open-source tool for measuring JavaScript code coverage. We report the average results over five runs to account for the non-determinism behaviour that stems from crawling the application. In addition, we assess each step in our approach separately as follows: (1) compare the statement coverage achieved by our function coverage maximization with a method that chooses the next

state/event for the expansion uniformly at random, (2) evaluate the effectiveness of applying mutation techniques (Algorithm 2) to reduce the number of assertions generated.

**Test Oracles (RQ2).** To evaluate the fault finding capability of JSEFT (RQ2), we simulate web application faults by automatically seeding each application with 50 random faults according to the following fault category:

1. Changing conditional statements by modifying the upper/lower bounds of loop statements, changing the condition itself, as well as swapping consecutive conditional statements;
2. Modifying the values of global/local variables, and removing or changing their names, as well as modifying arithmetic operations;
3. Changing function parameters or function call arguments by swapping, removing, and renaming parameters/arguments. Changing the sequence of function calls within a given function if applicable;
4. Modifying DOM related properties.

The first three categories target JavaScript code while the last one targets both JavaScript and HTML code-levels. We automatically pick a random program point and seed a fault at that point according to our fault category. While mutations used for oracle generation have been selectively generated (as discussed in Section 6), mutations used for the purpose of evaluation are randomly generated from the entire application. Note that if the mutation used for the purpose of evaluation and the mutation used for generating oracles happen to be the same, we remove the mutant from the evaluation set. Next we run the whole generated test suite (including both function-level and event-based test cases) on the faulty version of the application. The fault is considered detected if an assertion generated by JSEFT fails and our manual examination confirms that the failed assertion is detecting the seeded fault. We measure the precision and recall as follows:

**Precision** is the rate of injected faults found by the tool that are actual faults:  $\frac{TP}{TP+FP}$

**Recall** is the rate of actual injected faults that the tool finds:  $\frac{TP}{TP+FN}$

where  $TP$  (true positives),  $FP$  (false positives), and  $FN$  (false negatives) respectively represent the number of faults that are correctly detected, falsely reported, and missed.

**Function State Reduction (RQ3).** To assess the efficacy of the function state reduction method (Algorithm 1), we compare both the statement coverage and fault finding capability, before and after applying the state reduction technique. The former includes all possible function states converted to a test case, while the latter includes reduced test cases after applying the state reduction technique. To measure the fault finding capability, the same set of mutations (50) is used in both cases using the same fault seeding procedure as described in RQ2. Since the function state reduction method is used for generating function-level test cases, we run only the function-level test suite on the mutated versions for this research question.

**Comparison (RQ4).** To assess how JSEFT performs with respect to existing JavaScript test automation tools, we compare its coverage and fault finding capability to that of ARTEMIS [7]. Similar to JSEFT, we give ARTEMIS 10 minutes in total for each application; we observed no improvements in the results obtained from running ARTEMIS for longer periods of time. We run ARTEMIS from the command line by setting the iteration option to 100 and enabling the coverage priority strategy, as described in [7]. Similarly, JSCover is used to measure the coverage of ARTEMIS (over 5 runs). We use the output provided by ARTEMIS to determine if the seeded mutations are detected by the tool, by following the same procedure as described above for JSEFT.

#### 4.3. Results

**Test Case Generation (RQ1).** Figure 5 depicts the statement coverage achieved by JSEFT for each application. The results show that the test cases generated by JSEFT achieve a coverage of 68.4% on average, ranging from 41% (ID 12) up to 99% (ID 1). We investigated why JSEFT has low

Table II. Results showing the effects of our **function coverage maximization**, **function state reduction**, and **mutation-based oracle generation** algorithms.

App ID	St. Coverage		State Reduction			Oracles		
	Fun. cov. maximize (%)	Random exploration (%)	#Func.States w/o reduction	#Func.States with reduction	Func.State Reduction (%)	#Assertions w/o mutation-reduction	#Assertions with reduction only	#Assertions with mutation-reduction
1	99	80	447	42	91	5101	273	171
2	78	78	828	23	97	23212	128	91
3	90	66	422	22	95	3520	137	69
4	75	75	43	18	58	1232	131	105
5	49	45	534	31	94	150	124	103
6	78	75	797	38	95	1648	203	149
7	63	58	1653	58	96	198202	522	365
8	56	50	32	23	28	78	71	67
9	82	82	1509	51	97	65403	417	259
10	71	69	71	22	69	6584	136	94
11	56	54	1383	136	90	2530	549	334
12	41	38	1530	64	96	3521	253	193
13	51	47	1401	161	88	2481	638	361
AVG	68.4	62.8	-	-	84.1	-	-	-

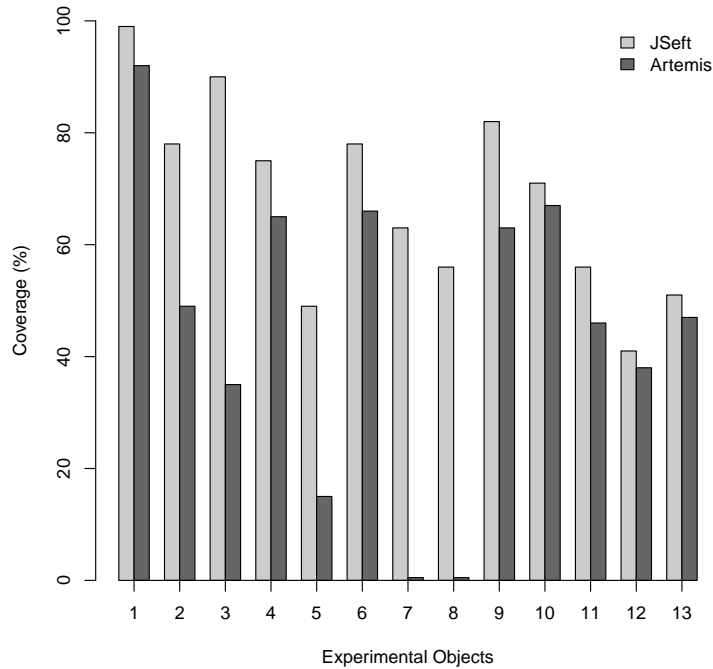


Figure 5. Statement coverage achieved.

coverage for some of the applications. For instance, we observed that in JointLondon (ID 8), the application contains JavaScript functions that are browser/device specific, i.e., they are exclusively

Table III. Fault detection.

App ID	# Injected Faults	JSEFT						JSEFT-all states			ARTEMIS	
		#FN	#FP	#TP	Precision (%)	Recall (%)	By func-level tests (%)	#FN	#TP	Recall (%)	Precision (%)	Recall (%)
1	50	0	0	50	100	100	30	0	50	100	100	20
2	50	7	0	43	100	86	73	7	43	86	100	12
3	50	4	0	46	100	92	17	4	46	92	100	8
4	50	12	0	38	100	76	28	11	39	78	100	22
5	50	24	0	26	100	52	25	24	26	52	100	0
6	50	9	0	41	100	82	15	9	41	82	100	16
7	50	16	0	34	100	68	24	14	36	72	100	0
8	50	21	0	29	100	58	26	21	29	58	100	0
9	50	6	0	44	100	88	41	6	44	88	100	24
10	50	16	0	34	100	68	65	13	37	74	100	8
11	50	19	0	31	100	62	27	19	31	62	100	6
12	50	25	0	25	100	50	17	25	25	50	100	22
13	50	23	0	27	100	54	26	20	30	60	100	28
AVG	-	14	0	36	100	72	32	13.3	36.7	73.4	100	12.8

executed in Internet Explorer, or iDevices. As a result, we are unable to cover them using JSEFT. We also noticed that some applications required more time to achieve higher statement coverage (e.g., in NarrowDesign ID 7), or they have a large DOM state space (e.g., BunnyHunt ID 5) and hence JSEFT is only able to cover a portion of these applications in the limited time it had available.

Table II columns under “St. Coverage” present JavaScript statement coverage achieved by our function coverage maximization algorithm versus a random strategy. The results show a 9% improvement on average, for our algorithm, across all the applications. We observed that our technique achieves the highest improvement when there are many dynamically generated clickable DOM elements in the application, for example, GhostBusters (ID 3). However, Tunnel (ID 2) and Fractal Viewer (ID 9) show no improvement in the statement coverage as these two applications have no dynamically generated or bound to event-listeners clickables. Instead, their few clickables are all placed in the HTML code of the application with a fixed event-handler per clickable. Thus, our approach achieves the same coverage as the random strategy for such applications.

**Fault finding capability (RQ2).** Table III presents the results on the fault finding capabilities of JSEFT. The table shows the total number of injected faults, the number of false negatives, false positives, true positives, and the precision and recall of JSEFT.

JSEFT achieves 100% precision, meaning that all the detected faults reported by JSEFT are real faults. *In other words, there are no false-positives.* This is because the assertions generated by JSEFT are all stable i.e., they do not change from one run to another. However, the recall of JSEFT is 72% on average, and ranges from 50 to 100%. This is due to false negatives, i.e., missed faults by JSEFT, which occur when the injected fault falls in either in the uncovered region of the application, or is not properly captured by the generated oracles.

The table also shows that on average 32% percent of the injected faults (ranges from 15–73%) are detected by function-level test cases, but not by our DOM event-based test cases. This shows that a considerable number of faults do not propagate to observable DOM states, and thus cannot be captured by DOM-level event-based tests. For example in the SimpleCart application (ID 10), if we mutate the mathematical operation that is responsible for computing the total amount of purchased items, the resulting error is not captured by event-based tests as the fault involves internal computations only. However, the fault is detected by a function-level test that directly checks the returned value of the function. This points to the importance of incorporating function-level tests in addition to event-based tests for JavaScript web applications. We also observed that even when an event-based test case detects a JavaScript fault, localizing the error to the corresponding JavaScript



code can be quite challenging. However, function-level tests pinpoint the corresponding function when an assertion fails, making it easier to localize the fault.

**Function State Reduction (RQ3).** The last three columns of Table II, under “Oracles”, present the number of assertions obtained by capturing the whole application’s state, (1) without applying any mutation-based oracle generation or using the function state reduction strategy (baseline), (2) with enabling the state reduction technique only, and (3) with applying both the state reduction and the mutation-based oracle generation algorithm, respectively. The results show that the number of assertions is decreased by 71.5% on average by using the state reduction technique compared to the baseline. The number of assertions is further decreased by 30% when applying the combination of mutation and state reduction strategy. This indicates that the state reduction technique plays a more prominent role than mutation in reducing the number of assertions. The reason is that the function state reduction technique removes considerable number of states each containing a number of oracles, while the mutation based method selectively chooses oracles within each state. We observe the most significant reduction of assertions for NarrowDesign (ID 7) from more than 198,000 to 365. NarrowDesign contains a number of functions that constantly perform display, hide, and resize activities. This results in creating considerable number of function states many of which are similar. Our state reduction eliminates such similar states, and hence results in significant reduction of assertions.

The columns under “State Reduction” in Table II show the number of function states before and after applying the function state reduction algorithm. The results show that the reduction strategy reduces the number of function states by 84.1% on average, with the maximum reduction being 97% in the case of Tunnel (ID 2) and FractalViewer (ID 9), and minimum reduction being 28% in the case of JointLondon (ID 8). As shown in Table II, JointLondon has the minimum number of function states as well as assertions even before performing any state reduction or mutation analysis. In addition to the application’s few number of variables, to achieve various functionality each function is called only a few times which contribute to the lower number of states/assertions compared to other applications.

We evaluate the effect of the state reduction strategy on both the statement coverage and fault-finding capability of JSEFT. We find that it does not reduce the statement coverage as it includes at least one function state from each of the covered branch sets as described in Section 3.2. For the fault-finding capability, we measured the false-negatives, true positives and recall of JSEFT without applying the state reduction. These values are shown in the columns under “JSEFT-all states” in Table III. As shown in the table, there is only a marginal decrease in the fault finding capability (less than 2% in recall) when using the state reduction strategy.

**Comparison (RQ4).** Figure 5 shows the code coverage achieved by both JSEFT and ARTEMIS on the experimental objects running for the same amount of time, i.e., 10 minutes. The test cases generated by JSEFT achieve 68.4% coverage on average (ranging from 41–99%), while those generated by ARTEMIS achieve only 44.8% coverage on average (ranging from 0–92%). Overall, the test cases generated by JSEFT achieve 53% more coverage than ARTEMIS, which points to the effectiveness of JSEFT in generating high coverage test cases. Further, as can be seen in the bar plot of Figure 5, for all the applications, the test cases generated by JSEFT achieve higher coverage than those generated by ARTEMIS. This increase was more than 226% in the case of Bunnyhunt (ID 5). For two of the applications, NarrowDesign (ID 7) and JointLondon (ID 8), ARTEMIS was not able to complete the testing task within the allocated time of ten minutes. Thus we let ARTEMIS run for an additional 10 minutes for these applications (i.e., 20 minutes in total). Even then, neither application completes under ARTEMIS.

Table III shows the precision and recall achieved by JSEFT and ARTEMIS. With respect to fault finding capability, unlike ARTEMIS that detects only generic faults such as runtime exceptions and W3C HTML validation errors, JSEFT is able to accurately distinguish faults at the code-level and DOM-level through the test oracles it generates. Both tools achieve 100% precision, however, JSEFT achieves five-fold higher recall (72% on average) compared with ARTEMIS (12.8% recall on average).

## 5. DISCUSSION

In this section, we discuss some of the limitations of our current implementation and threats to validity of our results.

### 5.1. Limitations

**Direct testability of functions.** Using our technique, it is not possible to produce test cases directly for all JavaScript functions. Non-testable JavaScript functions include (1) anonymous functions, (2) private function closures, and (3) functions that uses the `this` keyword as a reference to the object of which that function is a property/method. Depending on where the function is called from, the value of `this` can be different. While these functions can be called at the highest program scope they belong to, it is not possible to call them directly in test cases, which makes it challenging to assess their outcomes.<sup>†</sup> This is especially so if the tester is interested in stand-alone examination of these functions. Based on our observations during the evaluation, there is a non-trivial amount of non-testable JavaScript code out there. One possible future direction is to automatically refactor non-testable JavaScript code to make it testable.

**Undetected faults.** As mentioned in Section 4.3, our test cases miss some of the faults i.e., result in false negatives. BunnyHunt (ID 5) and TuduList (ID 12) achieve the lowest fault detection recall, which relates to their low coverage compared to the other applications. We describe possible reasons behind the observed low coverage for these applications, and possible ways to alleviate them.

1. *Unstable Behaviour of Test Cases:* Testing the correct behaviour of a function requires stable assertions, meaning that the test case and its assertions should always be deterministic. Thus, to achieve stability we need to remove, (1) random generation functions as well as functions that rely on the output of a randomized function, and (2) functions or event sequences that exhibit unstable behaviour in different executions due to highly dynamic nature of JavaScript. Traditional event based testing is not affected by such unstable behavioural characteristics of the functions as they run the events and compare the whole captured state after each event is triggered. One way to alleviate this would be to come up with invariants that characterize these functions and events, and use them for testing.
2. *Other limitations:* Our implementation is currently not able to trigger events that depend on the keyboard movements or drag and drop actions. Therefore, we obtain lower coverage for applications that contain such events due to the presence of functions that are either sensitive to keyboard movements or executed during drag and drop actions. This can be alleviated by a more comprehensive implementation - this is not a fundamental issue.

### 5.2. Threats to Validity

An external threat to the validity of our results is the limited number of web applications that we use to evaluate our approach. We mitigated this threat by using JavaScript applications that cover various application types and functionality. Another external threat is that we validate the failed assertions through manual inspection, which can be error-prone. To mitigate this threat, we carefully inspected the code in which the assertion failed to make sure that the injected fault was indeed responsible for the assertion failure.

An internal threat to validity is that we inject mutations to assess the fault finding capability of JSEFT, which can be biased towards JSEFT. However, we have attempted to mitigate this threat by following a standard methodology for fault seeding using a random mutation location selection method and a well-defined fault category. Moreover, during the evaluation of JSEFT, if the mutation used to evaluate the fault finding capability of the tool happens to be the same the mutation used to generate oracles, the mutant is removed from the evaluation set.

<sup>†</sup>This is similar to testing private methods in e.g., Java

Finally, regarding the reproducibility of our results, JSEFT and all the applications used in this study are publicly available, thus making the study replicable.

## 6. RELATED WORK

**Web application testing.** Andrews et al. [6] propose hierarchical Finite State Machine (FSM) to mitigate the state space explosion problem in web application testing, by reducing the number of states and transitions in the generated FSM. However, this technique is not fully automated as the construction of FSM requires manual effort. Marchetto and Tonella [20] propose a search-based algorithm for generating event-based sequences to test Ajax applications. To test PHP web applications Alshahwan and Harman [5] propose SWAT, a search-based method by using static analysis. Their main goal is to maximize the branch coverage. Heidegger et al. propose a random testing framework, called JSConTest [17]. The authors propose a contract language for JavaScript that allows a programmer to annotate the program with functional contracts. Annotations are used to guide the random test generation in finding counter examples. Praphamontripong et al. [30] apply mutation testing for Java Server Pages (JSP) and Java Servlets. The mutation operators are specifically proposed for these types of applications. Their results show that using mutation testing is particularly effective in detecting web application errors that occur on the server-side. Mesbah et al. [22] apply dynamic analysis to construct a model of the application's state space, from which event-based test cases are automatically generated. In subsequent work [23], they propose generic and application-specific invariants as a form of automated soft oracles for testing AJAX applications. Our earlier work, JSART [24], automatically infers program invariants from JavaScript execution traces and uses them as regression assertions in the code. Sen et al. [32] recently proposed a record and replay framework called Jalangi. It incorporates selective record-replay as well as shadow values and shadow execution to enable writing of heavy-weight dynamic analyses. The framework is able to track generic faults such as `null` and `undefined` values as well as type inconsistencies in JavaScript. Jensen et al. [18] propose a technique to test the correctness of communication patterns between client and server in AJAX applications by incorporating server interface descriptions. They construct server interface descriptions through an inference technique that can learn communication patterns from sample data. Li et al. [19] propose SymJS, a symbolic execution engine for JavaScript web applications. SymJS uses symbolic and dynamic feedback directed event space exploration as well as dynamic taint analysis to improve event sequence generation. Milani Fard et al. [11] propose a concolic-based tool, called ConFix to automatically generate DOM fixtures and DOM dependent function arguments for JavaScript unit tests. Saxena et al. [31] combine random test generation with the use of symbolic execution for systematically exploring a JavaScript application's event space as well as its value space, for security testing. Our work is different in that they either target the generation of event sequences at the DOM level, while we also generate unit tests at the JavaScript code level, which enables us to cover more and find more faults, or they do not address the problem of test oracle generation and only check against soft oracles (e.g., invalid HTML). In contrast, we generate strong oracles that capture application behaviours, and can detect a much wider range of faults.

Perhaps the most closely related work to ours is ARTEMIS [7], which supports automated testing of JavaScript applications. ARTEMIS considers the event-driven execution model of a JavaScript application for feedback-directed testing. In this paper, we quantitatively compare our approach with that of ARTEMIS (Section 4).

**Oracle generation.** There has been limited work on oracle generation for testing. Fraser et al. [15] propose  $\mu$ TEST, which employs a mutant-based oracle generation technique. It automatically generates unit tests for Java object-oriented classes by using a genetic algorithm to target mutations with high impact on the application's behaviour. They further identify [14] relevant pre-conditions on the test inputs and post-conditions on the outputs to ease human comprehension. The authors assume that the tester will manually correct the generated oracles. However, the results on the effectiveness of such approaches which rely on the "generate-and-fix" assumption to construct

test oracles are not conclusive [13]. Xie et al. explore test oracle generation for GUI systems [35]. Eclat [28] is used for automatically generating invariant-based oracles. Differential test case generation approaches [34, 10] are similar to mutation-based techniques in that they aim to generate test cases that show the difference between two versions of a program. However, mutation-based techniques such as ours, do not require two different versions of the application. Rather, the generated differences are in the form of controllable mutations that can be used to generate test cases capable of detecting regression faults in future versions of the program.

Pastore et al. [29] exploit crowd sourcing approach to check assertions. The developer produces tests and provides sufficient API documentation for the crowd such that crowd workers can determine the correctness of assertions. However, recruiting qualified crowd to generate test oracles can be quite challenging. Staats et al. [33] address the problem of selecting oracle data, which is formed as a subset of internal state variables as well as outputs for which the expected values are determined. They apply mutation testing to produce oracles and rank the inferred oracles in terms of their fault finding capability. This work is different from ours in that they merely focus on supporting the creation of test oracles by the programmer, rather than fully automating the process of test case generation. Further, (1) they do not target JavaScript; (2) in addition to the code-level mutation analysis, we propose DOM-related mutations to capture error-prone [27] dynamic interactions of JavaScript with the DOM. Milani Fard et al. [12] propose Testilizer which utilizes DOM-based test suite of the web application to explore alternative paths and consequently regenerate assertions for new detected states. Our work is different from this approach in that we generate unit level code-based test cases and assertions.

## 7. CONCLUSIONS

In this paper, we presented a technique to automatically generate test cases for JavaScript applications at two complementary levels: (1) individual JavaScript functions, (2) event sequences. Our technique is based on algorithms to maximize function coverage and minimize function states needed for efficient test generation. We also proposed a method for effectively generating test oracles along with the test cases, for detecting faults in JavaScript code as well as on the DOM tree. We implemented our approach in an open-source tool called JSEFT. We empirically evaluated JSEFT on 13 web applications. The results show that the generated tests by JSEFT achieve high coverage (68.4% on average), and that the injected faults can be detected with a high accuracy rate (recall 70%, precision 100%). We also find that our approach outperforms an existing JavaScript test automation framework in terms of coverage and fault detection capability. Compared with ARTEMIS, our tool achieves 53% more statement coverage as well as five-fold higher recall.

Our future work will include applying a more robust mutation analysis technique to generate non-fragile DOM-level test assertions. Another possible future direction is to identify the extent of hard-to-test functions (e.g.; private function closures) written by JavaScript developers through an empirical study. Automated code refactoring systems can be used to expose such functions to unit tests. JavaScript developers can also make use of the results of empirical study as a coding recommendation to make their future applications more testable and consequently more maintainable.

## REFERENCES

1. jQuery API. <http://api.jquery.com>. Accessed: 2014-06-30.
2. JSCover. <http://tntim96.github.io/JSCover/>. Accessed: 2014-01-30.
3. JSeft. <http://salt.ece.ubc.ca/software/jseft/>. Accessed: 2014-06-30.
4. Qunit. <http://qunitjs.com>. Accessed: 2014-06-30.
5. N. Alshahwan and M. Harman. Automated web application testing using search based software engineering. In *Proc. IEEE International Conference on Automated Software Engineering (ASE'11)*, pages 3–12, 2011.
6. A. Andrews, J. Offutt, and R. Alexander. Testing web applications by modeling with FSMs. *Software and Systems Modeling*, 4(3):326–345, July 2005.

7. S. Artzi, J. Dolby, S. Jensen, A. Møller, and F. Tip. A framework for automated testing of JavaScript web applications. In *Proc. 33rd International Conference on Software Engineering (ICSE'11)*, pages 571–580, 2011.
8. T. H. Cormen, C. Stein, R. L. Rivest, and C. E. Leiserson. *Introduction to Algorithms*. McGraw-Hill Higher Education, 2nd edition, 2001.
9. R. DeMillo, R. Lipton, and F. Sayward. Hints on test data selection: Help for the practicing programmer. *Computer*, 11(4):34–41, 1978.
10. S. Elbaum, H. N. Chin, M. B. Dwyer, and M. Jorde. Carving and replaying differential unit test cases from system test cases. *IEEE Transactions on Software Engineering (TSE)*, 35(1):29–45, 2009.
11. A. M. Fard and A. Mesbah. Generating fixtures for JavaScript unit testing. In *Proc. IEEE International Conference on Automated Software Engineering (ASE'15)*, 2015.
12. A. M. Fard, M. Mirzaaghaei, and A. Mesbah. Leveraging existing tests in automated test generation for web applications. In *Proc. IEEE International Conference on Automated Software Engineering (ASE'14)*, pages 67–78, 2014.
13. G. Fraser, M. Staats, P. McMinn, A. Arcuri, and F. Padberg. Does automated white-box test generation really help software testers? In *Proc. International Symposium on Software Testing and Analysis (ISSTA'13)*, pages 188–198. ACM, 2013.
14. G. Fraser and A. Zeller. Generating parameterized unit tests. In *Proc. International Symposium on Software Testing and Analysis (ISSTA'11)*, pages 364–374, 2011.
15. G. Fraser and A. Zeller. Mutation-driven generation of unit tests and oracles. *IEEE Transactions on Software Engineering (TSE)*, 38(2):278–292, 2012.
16. M. J. Harrold, R. Gupta, and M. L. Soffa. A methodology for controlling the size of a test suite. *ACM Transactions on Software Engineering and Methodology (TOSEM)*, 2(3):270–285, 1993.
17. P. Heidegger and P. Thiemann. Contract-driven testing of JavaScript code. In *Proc. International Conference on Objects, Models, Components, Patterns (TOOLS'10)*, pages 154–172, 2010.
18. C. Jensen, A. Møller, and Z. Su. Server interface descriptions for automated testing of JavaScript web applications. In *Proc. European Software Engineering Conference and ACM SIGSOFT International Symposium on Foundations of Software Engineering (ESEC/FSE'013)*. ACM, 2013.
19. G. Li, E. Andreasen, and I. Ghosh. SymJS: Automatic symbolic testing of JavaScript web applications. In *Proc. ACM SIGSOFT International Symposium on the Foundations of Software Engineering (FSE'14)*, pages 449–459, 2014.
20. A. Marchetto and P. Tonella. Using search-based algorithms for Ajax event sequence generation during testing. *Empirical Software Engineering*, 16(1):103–140, 2011.
21. A. Marchetto, P. Tonella, and F. Ricca. State-based testing of Ajax web applications. In *Proc. 1st Int. Conference on Sw. Testing Verification and Validation (ICST'08)*, pages 121–130. IEEE Computer Society, 2008.
22. 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)*, 6(1):3:1–3:30, 2012.
23. A. Mesbah, A. van Deursen, and D. Roest. Invariant-based automatic testing of modern web applications. *IEEE Transactions on Software Engineering (TSE)*, 38(1):35–53, 2012.
24. S. Mirshokraie and A. Mesbah. JSART: JavaScript assertion-based regression testing. In *Proc. International Conference on Web Engineering (ICWE'12)*, pages 238–252. Springer, 2012.
25. S. Mirshokraie, A. Mesbah, and K. Pattabiraman. Efficient JavaScript mutation testing. In *Proc. 6th International Conference on Software Testing Verification and Validation (ICST'13)*, pages 74–83. IEEE Computer Society, 2013.
26. S. Mirshokraie, A. Mesbah, and K. Pattabiraman. JSEFT: Automated javascript unit test generation. In *Proc. 8th International Conference on Software Testing Verification and Validation (ICST'15)*, pages 1–10. IEEE Computer Society, 2015.
27. F. Ocariza, K. Bajaj, K. Pattabiraman, and A. Mesbah. An empirical study of client-side JavaScript bugs. In *Proc. ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM'13)*, pages 55–64. IEEE Computer Society, 2013.
28. C. Pacheco and M. Ernst. Eclat: Automatic generation and classification of test inputs. In *Proc. European Conference on Object-Oriented Programming (ECOOP'05)*, pages 50–527, 2005.
29. F. Pastore, L. Mariani, and G. Fraser. CrowdOracles: Can the crowd solve the oracle problem? In *Proc. International Conference on Software Testing Verification and Validation (ICST'13)*, pages 342–351. IEEE Computer Society, 2013.
30. U. Praphamontipong and J. Offutt. Applying mutation testing to web applications. In *Proc. IEEE International Conference on Software Testing, Verification and Validation Workshops (ICSTW'10)*, pages 132–141, 2010.
31. P. Saxena, D. Akhawe, S. Hanna, F. Mao, S. McCamant, and D. Song. A symbolic execution framework for JavaScript. In *Proc. Symp. on Security and Privacy (SP'10)*, pages 513–528. IEEE Computer Society, 2010.
32. K. Sen, S. Kalasapur, T., and S. Gibbs. Jalangi: A selective record-replay and dynamic analysis framework for JavaScript. In *Proc. European Software Engineering Conference and ACM SIGSOFT International Symposium on Foundations of Software Engineering (ESEC/FSE'013)*. ACM, 2013.
33. M. Staats, G. Gay, and M. Heimdahl. Automated oracle creation support, or: How I learned to stop worrying about fault propagation and love mutation testing. In *Proc. International Conference on Software Engineering (ICSE'11)*, pages 870–880, 2011.
34. K. Taneja and T. Xie. DiffGen: Automated regression unit-test generation. In *Proc. International Conference on Automated Software Engineering (ASE'08)*, pages 407–410. IEEE Computer Society, 2008.
35. Q. Xie and A. M. Memon. Designing and comparing automated test oracles for gui-based software applications. *ACM Transactions on Software Engineering and Methodology (TOSEM)*, 16(1), 2007.

## 8. EXTENSIONS W.R.T. THE ICST'15 PAPER

(This appendix is for reviewers only, and will not be included in the final version).

This paper is a revised and extended version of our paper 'JSEFT: Automated JavaScript Unit Test Generation', which appeared in the *Proceedings of the International Conference on Software Testing, Verification and Validation (ICST 2015)* [26].

The main extensions include:

- Improving the approach as follows:
  1. Revising the function state reduction technique (explained in Section 3.3.1); The approach diagram (Figure 2), and the function state reduction algorithm (Algorithm 1) are revised to reflect changes to the state reduction technique. In accordance with the revised approach diagram, we restructured the approach (Section 3) into three sub-sections: (1) maximizing function coverage, (2) extracting DOM and JavaScript function states, and (3) mutation analysis.
- Extending the evaluation section as follows:
  1. Comparing the fault finding capability obtained by (1) using the state reduction mechanism and (2) including the whole application's state.
  2. Comparing the number of assertions in the following scenarios: (1) capturing the whole application's state, (2) with enabling the state reduction technique only, and (3) with applying the state reduction as well as the mutation-based oracle generation algorithm;
- Adding one more research question to the evaluation section (RQ3 in Section 4) to discuss the results of the new comparisons.
- The new results are presented in Section 4.3. The comparative results are presented in Table II and Table III.
- Adding a discussion section (Section 5).
- A more elaborate coverage of related work (Section 6).
- Furthermore, we made a full pass over the text, leading to many improvements throughout the paper, including various changes that originated from discussions at the ICST'15 conference and other venues where this work was presented.

While difficult to measure objectively, we believe this accounts for an extension of over 30% – a common guideline used for journal versions of papers that are based on conference papers.