

# Generating Event Sequence-Based Test Cases Using GUI Run-Time State Feedback

Xun Yuan and Atif M Memon

Department of Computer Science, University of Maryland,

College Park, MD 20742, USA

{xyuan, atif}@cs.umd.edu

## Abstract

Graphical user interfaces (GUIs) are the sole mode of interaction between end-users and back-end code for almost all of today’s software applications. Because of this strategic role of GUIs, their quality has become important. During GUI testing, test cases—modeled as sequences of user events—sample the vast input space of all possible sequences with the goal of detecting faults; for effective testing, it is important to sample this space carefully. Existing techniques either sample manually or employ manually constructed abstract models—the abstraction and/or subsequent test-case generation algorithms enable sampling.

This paper presents an alternative approach to GUI testing – it’s focus is on developing a *fully automatic* model-driven technique to generate GUI test cases. The technique is novel in that it uses feedback from the execution of a “*seed test suite*” on a GUI. The seed suite is generated automatically using an existing structural *event-interaction graph* (EIG) model of the GUI. During its execution, the run-time effect of each GUI event on all other events pinpoints new important *event-semantic interaction* (ESI) relationships between them, which are used to automatically create an *ESI graph* (ESIG) model and generate new test cases. Together with a reverse-engineering algorithm used to obtain the EIG, seed suite, ESIG, and new test cases, the feedback-based technique yields a fully automatic, end-to-end GUI testing process. Two independent studies on eight applications demonstrate that the feedback-based technique (1) is able to significantly improve existing techniques and help identify serious problems in the software and (2) the ESI relationships captured via GUI state yield test suites that most

1 often detect more faults than their code-, event-, and event-interaction-coverage equivalent  
2 counterparts.

## 3 **1 Introduction**

4 Automated test case generation (ATCG) has become increasingly popular due to its potential to  
5 reduce testing cost and improve software quality [22]. A typical approach used for ATCG is  
6 to create an abstract model (*e.g.*, state-machine model [8, 67, 73], event-flow model [44]) of the  
7 application under test (AUT) and employ the model to generate test cases. While successful at  
8 reducing overall testing cost, in practice, ATCG continues to be resource-intensive, especially to  
9 create and maintain the model. A few researchers have recognized that the tasks of model creation  
10 and maintenance may be aided by leveraging the execution results of some existing test cases.  
11 Consequently, they have developed *automated feedback-based techniques* to augment the models  
12 [14, 21, 25–27, 37, 52, 54, 83, 84, 86]. These techniques require an initial test case/suite to be created,  
13 either manually or automatically, and executed on the software. Feedback from this execution is  
14 used to *automatically* generate additional test cases. The nature of feedback depends largely on the  
15 goal of the ATCG algorithm. A common example of feedback is a code coverage report used to  
16 automatically generate additional test cases that improve overall test coverage [25–27, 37, 52, 54].  
17 Few techniques use feedback from the AUT’s *run-time state* to generate additional test cases, *e.g.*,  
18 in the form of outcomes of programmer-supplied predicates in the code to cover all non-isomorphic  
19 inputs [14], operational abstractions to cover increased program behaviors [21, 84], and partially  
20 generated non-exception-throwing method-call sequences to generate longer sequences [57].

21 This paper presents a new feedback-based technique for automated testing of graphical user  
22 interfaces (GUIs). The technique starts with an automatically reversed engineered *structural* model  
23 of the GUI, employs the model to automatically generate specific types of test cases (sequences of  
24 GUI events that exercise GUI widgets), executes them, and uses the execution results as feedback  
25 to automatically generate additional test cases. The feedback is an abstraction of the run-time  
26 state of GUI widgets. As noted by Mathur [41], there is a strong relationship between events in a

1 software and its states; GUIs are no exception. Because GUI events may drive the software into  
2 potentially new states, this feedback-based approach leverages run-time GUI state information to  
3 prune the state-space, thereby helping to reduce the number of test cases that need to be executed.

4 The nature of GUIs, their test cases, and the maturity of our existing model-based GUI test-  
5 case generation algorithms lend themselves to feedback-based techniques for a number of reasons.  
6 First, GUI testing is extremely important because GUIs are used as front-ends to most software  
7 applications and constitute as much as half of software’s code [56]. A correct GUI is necessary  
8 for trouble-free execution of the application’s underlying “business logic” [10, 67, 73]. Second,  
9 our existing fully automatic model-based GUI test-case generation algorithms produce test cases  
10 that exhaustively test *two-way interactions* between GUI events; these test cases are called smoke  
11 tests [45]. They are used as the basis for feedback collection, *i.e.*, they form the seed suite. Finally,  
12 our existing tools are easily adapted to monitor and store the run-time state of the GUI.

13 Our previous empirical studies showed that smoke test cases reveal a large number of GUI  
14 faults; we and other researchers have shown that additional faults may be detected by testing  
15 certain types of multi-way interactions [10, 49]. The challenge, of course, is how to systematically  
16 generate test cases for these interactions. Exhaustively testing them is impossible because the  
17 number of GUI test cases grows exponentially with number of events in the test case. A practical  
18 alternative is to identify small subsets of events that interact in interesting ways with one another  
19 and hence should be tested together, and generate test cases that test multi-way interactions among  
20 members of each subset. In this paper, we use the feedback-based technique to *automatically*  
21 identify such sets of events.

22 The new feedback-based technique has been used in a fully automatic end-to-end process for  
23 a specific type of GUI testing. The seed test suite (in this case the smoke tests) is generated  
24 automatically using an existing *event-interaction graph* (EIG) model of the GUI. The EIG is a  
25 structural model of the GUI. More specifically, it represents *all possible sequences* of events that  
26 may be executed on the GUI. Note that an EIG has the flavor of the conventional control-flow  
27 model that represents all possible execution paths in a program [3], and a data-flow model that

1 represents all possible definitions and uses of a memory location [62]; except that an EIG is at the  
2 GUI event level of abstraction.

3 The EIG-generated smoke suite is executed on the GUI using an automatic test case replayer.  
4 During test execution, the run-time state of GUI widgets is collected and used to automatically  
5 identify an *Event Semantic Interaction* (ESI) relationship between pairs of events. This relationship  
6 captures how a GUI event is related to another in terms of how it modifies the other’s execution  
7 behavior. Informally, event  $e_x$  is ESI-related to  $e_y$  iff  $e_x$  influences the run-time behavior of  $e_y$ ,  
8 where “run-time behavior” is evaluated in terms of properties of GUI widgets.

9 The ESI relationships are used to automatically construct a new model of the GUI, called the  
10 *Event Semantic Interaction Graph* (ESIG). Because the seed suite is generated from the EIG (a  
11 structural model) and the ESI relationship is obtained in terms of event execution (a dynamic  
12 activity), the ESIG captures certain structural and dynamic aspects of the GUI. The ESIG is then  
13 used to automatically generate new test cases. These test cases have an important property – each  
14 event is ESI-related to its subsequent event, *i.e.*, it was shown to influence the subsequent event  
15 during execution of the seed suite.

16 This entire process, including the scripts required to set up, execute, and tear down test cases,  
17 has been implemented and executes without human intervention. Two independent studies have  
18 been conducted on eight GUI-based Java applications to evaluate and understand this new ap-  
19 proach.

20 In an earlier report of this work [86], we described the first study, which used four well-tested  
21 and popular applications downloaded from SourceForge; the study demonstrated that the feedback-  
22 based technique improves our existing techniques with little additional cost. The ESI relationship  
23 is successful at identifying complex interactions between GUI event handlers that lead to serious  
24 failures. We presented details of some failures, emphasizing on why they were not detected by  
25 the earlier techniques. The failures were reported on the SourceForge bug reporting site;<sup>1</sup> in re-  
26 sponse, the developers fixed some of the bugs. The developers had never detected our reported

---

<sup>1</sup>For example, [https://sourceforge.net/tracker/?func=detail&atid=535427&aid=1536078&group\\_id=72728](https://sourceforge.net/tracker/?func=detail&atid=535427&aid=1536078&group_id=72728).

1 failures before because their own tools and testing processes were unable to comprehensively and  
2 automatically test the applications.

3 We now extend the research with the second study, conducted on four fault-seeded Java appli-  
4 cations developed in house; this study shows that (1) the automatically identified ESI relationships  
5 between events help to generate test suites that detect more faults than their code-, event-, and  
6 event-interaction-coverage equivalent counterparts, (2) certain characteristics of the seeded faults  
7 prevent their detection by the earlier technique, but not the new technique, (3) several of our  
8 missed faults remain undetected because of limitations with our automated GUI-based test oracle  
9 (a mechanism that determines whether a test case passed or failed), and (4) several of the remaining  
10 undetected faults require long event sequences.

11 Finally, we note that the use of software models to generate sequences of events (commands,  
12 method calls, data inputs) for software testing is not new. Numerous researchers have developed  
13 techniques that employ state machine models [1, 11, 17–20, 23, 68, 74], grammars [5, 31, 42, 71, 72],  
14 AI planning [30, 39, 47, 66], genetic algorithms [35], probabilistic models [75, 77–79], architecture  
15 diagrams [65], and specifications [32, 33] to generate such sequences. All the above techniques are  
16 useful, in that they can be used to generate different types of test cases for different domains. All of  
17 them are based on manually created models. The research presented in this paper is orthogonal to  
18 the other model-based techniques; we focus on enhancing an existing model (in our case the model  
19 is obtained automatically) via test execution feedback. In particular, we leverage our existing  
20 graph-traversal techniques based on an automatically obtained reverse engineered GUI model [46,  
21 50] to develop a fully automatic model-based testing technique. Run-time feedback is used to  
22 enhance the model and generate new test cases. We feel that this type of approach may be used  
23 for the other model-based techniques mentioned above – these other models may also be enhanced  
24 with software execution and test execution feedback.

25 The main contributions of this work include:

- 26 • extension of our previous work on automated, model-based, systematic GUI test-case gen-  
27 eration,

- 1 • definition of new relationships among GUI events based on the GUI widgets that they  
2 use/influence,
- 3 • utilization of run-time state to explore a larger input space and improve fault-detection ef-  
4 fectiveness,
- 5 • immersion of the feedback-based technique into a fully automatic end-to-end GUI testing  
6 process and demonstration of its effectiveness on fielded and fault-seeded applications,
- 7 • empirical evidence tying fault characteristics to the type of test suites, and
- 8 • demonstration that certain faults require well-crafted combinations of test cases and test  
9 oracles.

10 The next section discusses related literature. Section 3 introduces basic GUI concepts and  
11 reviews the EIG model that forms the basis of the new ESIG model. Section 4 defines the ESI  
12 relationship and uses it to define an ESIG. Sections 5 and 6 evaluate the new feedback-based  
13 technique. Finally, Section 7 concludes with a discussion of future work.

## 14 **2 Related Work**

15 To the best of our knowledge, this is the first work that utilizes run-time information as feedback  
16 for model-based GUI test-case generation. However, run-time information has previously been  
17 employed for various aspects of test automation, and model-based testing has been applied to  
18 conventional software as well as *event-driven software* (EDS). This section presents an overview  
19 of related research in the areas of model-based and EDS testing, GUI testing, and the use of run-  
20 time information as feedback for test generation.

### 21 **2.1 Model-based & EDS Testing**

22 Model-based testing automates some aspect of software testing by employing a model of the soft-  
23 ware. The model is an abstraction of the software's behavior from a particular perspective (*e.g.*,  
24 software states, configuration, values of variables, etc.); it may be at different levels of abstraction,

1 such as abstract states, GUI states, internal variable states, or path predicates. Models may be de-  
2 rived from a formal specification of the software or reverse engineered by observing the software's  
3 execution behavior. They may be described using various languages and mathematical objects.

4 **State Machine Models:** The most popular models used for software testing are *state machine*  
5 *models*. They model the software's behavior in terms of its abstract or concrete states; they are  
6 typically represented as state-transition diagrams. Several types of state machine models have  
7 been used for software testing, such as *Finite State Machine Models* (FSM) [4, 6, 24, 28], *UML*  
8 *Diagram-based Models* [40] and *Markov Chains* [34, 76].

9     There have been numerous reports of success stories using FSM models for test automa-  
10 tion. For example, Microsoft researchers [6] modeled the control flow of an object-oriented soft-  
11 ware under test as an FSM and described it using the *Abstract State Machine Language* (AsmL  
12 `research.microsoft.com/fse/asml`). A traversal engine (part of Spec Explorer, a tool  
13 for advanced model-based specification and conformance testing), used the resulting finite state  
14 machine to produce behavioral tests to cover all explored transitions. Hong *et al.* also used FSMs  
15 for unit testing of classes in object-oriented programs; they used FSMs to model interactions be-  
16 tween class data members and member functions [28]. The FSMs, called *class state machines*  
17 (CSM), were then transformed into *class flow graphs* (CFG); test case generation was done by  
18 selecting test cases according to the locations of definitions and uses of variables in the CFG. In  
19 another reported research, Farchi *et al.* used FSM models to test implementations of the *POSIX*  
20 standard and *Java* exception handling [24]. Both state machine models were created from the  
21 software specifications and represented using the *GOTCHA Definition Language*. The GOTCHA-  
22 TCBean test generator was then used to automatically explore the state space from the model and  
23 generate an abstract test suite.

24     Various extensions of FSMs have also been used for testing. These extensions use variables  
25 to represent *context* in addition to states; the goal is to reduce the total number of states by using  
26 an orthogonal mechanism, in the form of explicit variables, to select state transitions. For exam-  
27 ple, an *extended finite state machine* (EFSM) makes use of a data state along with the input for

1 state transformation [4]; this EFSM is used by a tool called TestMaster to generate test cases by  
2 traversing all paths from the start state to the exit state.

3 Because test cases for EDS are sequences of events, many practitioners and researchers have  
4 found it natural to use state machine models for testing EDS [16,36,38,53]. The EDS is modeled in  
5 terms of states; events form transitions between states. Algorithms traverse these machine models  
6 to generate sequences of events. For example, Campbell *et al.* have applied state machine models  
7 to test object-oriented reactive systems [16]. Object states are modeled in terms of instance variable  
8 values; transitions are obtained from method invocations; test cases are sequences of method calls  
9 and are generated by traversing the model.

10 **Table-based Models:** Table-based models define software behavior in the form of tables relating  
11 model elements such as system modes, conditions, events and terms. These tables are then used as  
12 the basis for test case generation. The table-based modeling approach SCR, which is an abbrevia-  
13 tion for *software cost reduction*, has been used for security functional testing [12] and *Mars Polar*  
14 *Lander* software to detect faults [13].

15 **Grammars:** Production grammars have been used to test large, complex and safety-critical soft-  
16 ware systems; a popular example is the Java Virtual Machine [69]. The grammars are collections  
17 of non-terminal to terminal mappings that resemble regular parsing grammars. A production gram-  
18 mar produces a program (*i.e.*, a set of terminals, or tokens) starting from a high-level description  
19 (*i.e.*, a set of non-terminals). The composition of the generated programs models the restrictions  
20 placed on the software by the production grammar.

21 **Summary:** The above model-based testing techniques rely heavily on the manual or semi-manual  
22 construction of the abstract model. Consequently, they are prone to errors. Moreover, any change  
23 to the software requires reconstruction of the model, which is typically cumbersome and time  
24 consuming.



## 2.2 GUI Test Case Generation

Several automated techniques have been developed for GUI test case generation. All of them use a model of the software and algorithms to generate test cases from the model.

**State-Based Techniques:** Finite state machines have been used to model GUIs [7, 73]. A GUI's state is represented in terms of its windows and widgets; each user event triggers a transition in the FSM. For the sake of test-case generation, a test case is a sequence of user events and corresponds to a path in the FSM. As is the case for conventional software (discussed in the previous section), FSMs for GUIs also have scaling problems; this is due to the large number of possible states and user events in modern GUIs. Several GUI-domain-specific attempts have been made to handle the scalability issue. For example, Belli [7] converted a GUI FSM into simplified regular expressions. The regular expressions were used to generate event sequences. Shehady *et al.* [67] proposed variable finite state machine (VFSM), which augmented an FSM for a GUI with global variables that can assume a finite number of values during the execution of a test case. The value of each variable is used to determine the next state and output in response to an event. Event transition may modify values of these variables.

AI planning has also been used to manage the state-space explosion by eliminating the need for explicit states. AI planning models the infinite state space of a GUI [48]. A description of the GUI is manually created by a tester; this description is in the form of *planning operators*, which model the preconditions and effects (post-conditions) of each GUI event. Test cases are automatically generated from tasks (pairs of initial and goal states) by invoking a planner which searches for a path from the initial state to the goal state. However, the quality of the test cases is determined by the choice of tasks. Moreover, the manual operator definition and task selection may be expensive for large GUIs.

**Genetic Algorithm:** Test cases have been generated using genetic algorithms to mimic novice users [35]. The approach uses an expert to generate an initial event sequence manually and then uses genetic algorithm techniques to generate longer sequences. The assumption is that experts take a direct path when performing a task via the GUI, whereas novice users take longer, indirect

1 paths. Although useful for generating multiple test cases, the technique relies on an expert to gen-  
2 erate the initial test case. The final test suite depends largely on the paths taken by the expert user.  
3 The idea of using a task and generating an initial test case may be better handled by using plan-  
4 ning, since multiple test cases may be generated automatically according to some predetermined  
5 coverage criterion.

6 **Directed Graph Models:** In order to reduce manual work, several new systematic techniques  
7 based on graph models of the GUI have recently been developed. They are based on *Event Flow*  
8 *Graphs* (EFG) [45] and *Event Interaction Graphs* (EIG). Because of their central role in this paper,  
9 we discuss these models in Section 3.

10 **Summary:** FSM models and genetic algorithm based GUI testing suffer from the problem of  
11 manual creation of the model, *i.e.*, state machine and fitness function. The primary problem with  
12 the GUI graph models is that the number of event sequences grows exponentially with length.  
13 Hence, the existing graph-model based GUI test-case generation algorithms have only been able  
14 to generate test cases that cover all edges in the graph models, *i.e.*, they test *two-way* interactions  
15 between GUI events.

## 16 **2.3 Execution Feedback for Test Case Generation**

17 In an earlier report of this research [86], we introduced the idea of employing *feedback* from the  
18 execution of a seed test suite (our smoke tests generated using the EIG) to generate additional  
19 multi-way interaction test cases. A study on four large fielded open-source software applications  
20 demonstrated the feasibility and usefulness of this new approach; we showed that feedback was  
21 able to significantly improve our existing techniques and help identify/report serious problems in  
22 the software applications.

23 We now discuss similar feedback-based approaches used by other researchers on non-GUI  
24 software to generate test cases. Execution feedback refers to information obtained during test  
25 execution, and used to guide automatic test case/test suite generation. This is called *dynamic*  
26 *test case generation* and, to the best of our knowledge, was originally proposed by Miller and

1 Spooner [54]. In their technique, the software source code is instrumented to obtain execution  
2 feedback. The overall test case generation process starts by executing an initial test, which may  
3 be a test suite or a single test case. The execution feedback is collected and analyzed. The results  
4 are used to evaluate the “closeness,” according to some criterion, of the previous execution to the  
5 desired outcome; the model used to generate test cases is then modified accordingly and a new test  
6 case is generated. This loop stops when the “closeness” evaluation is satisfied.

7 Since then, several researchers have used similar ideas of dynamic test generation.

8 **Object Properties:** The work of Xie *et al.* is most closely related to this research [83, 84]. They  
9 have developed a framework that uses feedback in the form of *operational abstractions* (summaries  
10 of program run-time state) and object states to generate new test cases. This framework integrates  
11 specification-based test generation and dynamic specification inferences for test case generation.  
12 Specification-based test generation is based on formal specifications, which express the desired  
13 behavior of a program. However, because formal specifications are difficult to obtain, dynamic  
14 specification inference attempts to infer specifications, in the form of operational abstractions,  
15 automatically from software execution. The test case generation process starts from an existing  
16 test suite. Through executions of these test cases, object states (values of variables and parameters,  
17 and return values) are recorded at the entry and exit of method executions. Based on the collected  
18 traces and a set of pre-defined axiom-pattern templates, equality patterns are searched to create  
19 operational abstractions. The discovered operational abstractions consist of object properties that  
20 hold for all the observed executions. By removing or relaxing inferred preconditions on parameter  
21 values in the operational abstractions, both legal and illegal test cases are generated. The newly  
22 generated test cases are executed. Because they were generated by relaxing inferred preconditions,  
23 some of these test cases may cause an uncaught runtime exception. The other, non-crashing test  
24 cases are used to obtain new operational abstractions, which are again used to generate additional  
25 test cases. Other researchers have also used operational abstractions, combined with symbolic  
26 execution, to guide the generation of test cases [21].

27 **Method-call Sequences:** Pacheco *et al.* [58] have improved random unit test generation by

1 incorporating feedback obtained from executing test inputs as they are created. They build inputs  
2 incrementally by randomly selecting a method call to apply and finding arguments from among  
3 previously-constructed inputs. The key idea of their work is that they build upon a legal sequence  
4 of method calls, each of whose intermediate objects is sensible and none of whose methods throw  
5 an exception. As soon as an input is built, it is executed and checked against a set of contracts  
6 and filters. The result of the execution determines whether the input is redundant, illegal, contract-  
7 violating, or useful for generating more inputs. The technique outputs a test suite consisting of unit  
8 tests for the classes under test. Passing tests can be used to ensure that code contracts are preserved  
9 across program changes; failing tests (that violate one or more contract) point to potential errors  
10 that should be corrected.

11 Similarly, Boyapati *et al.* employ a feedback-based technique to obtain all non-isomorphic  
12 inputs (test cases) for a method [14]. A programmer develops (1) a “guided test generation engine”  
13 that outputs test cases to explore the method’s input space and (2) a predicate from the method’s  
14 preconditions to check the validity of the generated input. This technique prunes a large portion  
15 of the input space by monitoring the execution of the predicate on an initial test suite, guiding the  
16 engine and yielding a suite of all non-isomorphic inputs.

17 **Code Coverage Reports:** All other techniques in this category instrument elements (lines,  
18 branches, etc.) of the program code, execute an initial test case/suite, obtain a coverage report  
19 that contains the outcomes of conditional statements, and use automated techniques to generate  
20 better test cases. The techniques differ in their goals (*e.g.*, cover a specific program path, satisfy  
21 condition-decision coverage, cover a specific statement) and their test-case generation algorithms.  
22 For example, Miller *et al.* [54] use code coverage and decision outcomes to generate floating-point  
23 test data.

24 Several *iterative techniques* have been used to generate a test case that executes a given program  
25 path [26, 27, 37]. The generation is formulated as a function minimization problem. The gradient-  
26 descent approach is used to gradually adjust an initial test case so that it executes the given path.  
27 Control-flow information in the form of branch-predicate outcomes is collected during software

1 execution. One disadvantage of these approaches is that they can get stuck in a local minima during  
2 test case generation.

3 The *chaining* approach [25] has been used to generate test cases, each to cover a given pro-  
4 gram statement. An initial test case is executed; the program's control- and data-flow are used  
5 to determine whether the test case will lead to the given statement. If not, the branch function  
6 of the problematic branch is used to modify the test case. This process continues until the given  
7 statement is executed.

8 *Genetic algorithms* have also been used to automatically generate test suites that satisfy the  
9 *condition-decision* adequacy criterion [52], which requires that each condition in the program be  
10 true for at least one test case and false for at least one test case. A fitness function is defined  
11 for each branch. An initial test suite is obtained and executed. The fitness functions are used to  
12 evaluate the "goodness" of each test case. If a test case covers a new condition-decision, it is  
13 considered to be "more fit." The test cases in the gene pool evolve to obtain a new generation of  
14 test cases. The process stops until a desired level of fitness is obtained.

15 **Summary:** All the above execution feedback-based techniques have been used for a specific  
16 type of test case, that is, numerical data values. The feedback (in the form of branch predicate  
17 evaluations, condition-decision coverage, and object states) is used to tweak these numerical values  
18 in order to improve overall coverage. These techniques are not directly applicable to GUI testing  
19 because a GUI test case is a sequence of events. There is no clear notion of tweaking a GUI test  
20 case to improve coverage.

21 Although the techniques discussed in this section are not directly applicable to feedback-  
22 directed GUI test case generation, many of the underlying concepts have been used. For example,  
23 execution feedback is used to generate GUI test cases, the EIG model is used to generate the orig-  
24 inal seed suite, and traversal techniques from model-based testing are used to cover nodes and  
25 edges in the ESIG.

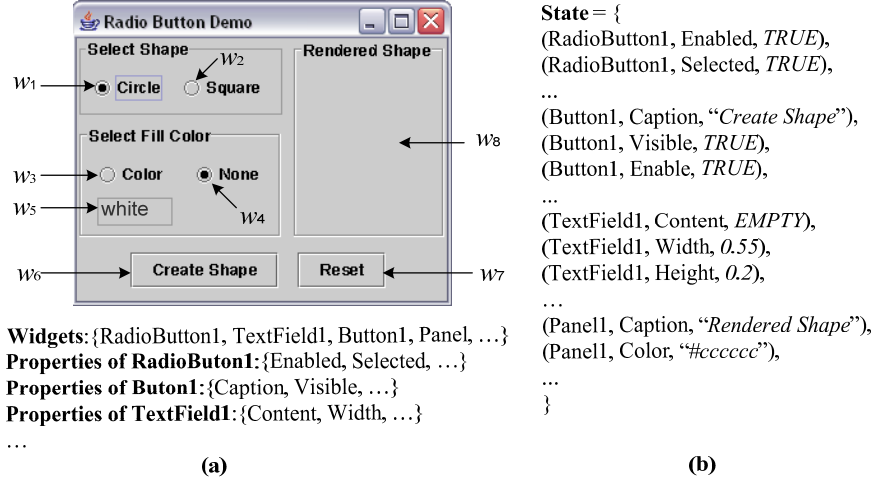


Figure 1: (a) "Radio Button Demo" GUI, (b) its Partial Observable State

### 3 Preliminaries

The feedback-based technique utilizes an abstraction of the GUI's *run-time state* collected and analyzed during the execution of test cases that cover *two-way interactions* between GUI events in order to generate test cases that test *multi-way interactions*. This section defines these terms and introduces notations for subsequent sections.

This work focuses on the class of GUIs that accept discrete events performed by a single user; the events are deterministic, *i.e.*, their outcomes are completely predictable.<sup>2</sup> A GUI in this class is composed of a set  $W$  of *widgets* (*e.g.*, buttons, text fields); each widget  $w \in W$  has a set  $P_w$  of *properties* (*e.g.*, color, size, font). At any time instant, each property  $p \in P_w$  has a unique *value* (*e.g.*, red, bold, 16pt); each value is evaluated using a function from the set of the widget's properties to the set of values  $V_p$ . The *GUI state* at any time instant is a set of triples  $(w, p, v)$ , where  $w \in W, p \in P_w$  and  $v \in V_p$ , *i.e.*, the observable state of the GUI. Figure 1 shows the partial GUI state of a simple application's window called Radio Button Demo. The GUI contains eight widgets labeled  $w_1$  through  $w_8$ ; a user can perform events  $e_1$  through  $e_7$  on  $w_1$  through  $w_7$ , respectively; no event can be performed on  $w_8$ .

A set of states  $S_I$  is called the *valid initial state set* for a particular GUI if the GUI may be in

<sup>2</sup>Testing GUIs that react to temporal and non-deterministic events and those generated by other applications is beyond the scope of this research.

1 any state  $S_i \in S_I$  when it is first invoked. The single *initial state* for `Radio Button Demo` has  
 2 `Circle` and `None` selected; the text-field corresponding to  $w_5$  has a default color “white”; and  
 3 the `Rendered Shape` area (widget  $w_8$ ) is empty.

4 The state of a GUI is not static; events  $e_1, e_2, \dots, e_n$  performed on the GUI change its state  
 5 and hence are modeled as functions that transform one state of the GUI to another. For `Radio`  
 6 `Button Demo`, event  $e_1$  sets the shape to a circle; if there is already a square in the `Rendered`  
 7 `Shape` area, then it is immediately changed to a circle. Event  $e_2$  is similar to  $e_1$ , except that it  
 8 changes the shape to a square. Event  $e_3$  enables the text-field  $w_5$ , allowing the user to enter a  
 9 custom fill-color, which is immediately reflected in the shape being displayed (if there is a shape  
 10 there). Event  $e_4$  reverts back to the “no fill color” state. Event  $e_5$  is used to fill a custom color in  
 11 the text-field  $w_5$ . Event  $e_6$  creates a shape in the `Rendered Shape` area according to current  
 12 settings of  $w_1 \dots w_5$ ; event  $e_7$  resets the entire software to its initial state.

13 GUIs contain two types of windows: (1) *modal windows*<sup>3</sup> (e.g., `FileOpen`, `Print`) that,  
 14 once invoked, monopolize the GUI interaction, restricting the focus of the user to the range of  
 15 events within the window until explicitly terminated (e.g., using `Ok`, `Cancel`), and (2) *modeless*  
 16 *windows* (e.g., `Find/Replace`) that do not restrict the user’s focus. If, during an execution of  
 17 the GUI, modal window  $\mathcal{M}_x$  is used to open another modal window  $\mathcal{M}_y$ , then  $\mathcal{M}_x$  is called the  
 18 *parent* of  $\mathcal{M}_y$  for that execution.

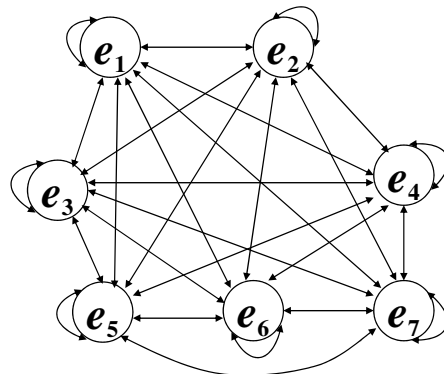


Figure 2: EIG of “Radio Button Demo” GUI

<sup>3</sup>Standard GUI terminology, e.g., see <http://java.sun.com/products/jlf/ed2/book/HIG.Dialogs.html>.

1 The seed test suite is generated using an *event-interaction graph* (EIG) model of the GUI,  
2 which is obtained automatically using a standard GUI-reverse-engineering algorithm [49]. The  
3 EIG abstraction of the GUI represents only two types of GUI events: *termination* and *system-*  
4 *interaction* events. Termination events close modal windows. Other *structural* events are used to  
5 open and close menus and modeless windows, and open modal windows, but are not represented  
6 in the EIG (for reasons presented in earlier work [49]). The remaining events, called system-  
7 interaction events, do not manipulate the structure of the GUI. Directed edges between nodes  
8 encode *execution paths*, *i.e.*, sequences of events, in the GUI. For example, an edge  $(e_x, e_y)$  shows  
9 that  $e_y$  may be executed after  $e_x$  along some *execution path*. The EIG for the `Radio Button`  
10 Demo is shown in Figure 2. Because this is a single-window GUI with no menus, it has no  
11 structural events; it contains only system-interaction events. There is one node for each of the  
12 widgets on which a user can perform an event.

13 The basic motivation behind using a graph model to represent a GUI is that various types  
14 of existing graph-traversal algorithms (with well-known run-time complexities) may be used to  
15 “walk” the graph, enumerating the events along the visited nodes, thereby generating test cases.  
16 In earlier research [49], an algorithm called `GenTestCases` was implemented that returned all  
17 possible paths (sequences of events) in the graph bounded to a specific length (number of EIG  
18 events) of 2. These length-2 sequences are said to test all *two-way interactions* between the EIG  
19 events. For the EIG of Figure 2, there are a total of 49 test cases of length 2, corresponding to the  
20 49 edges in the EIG. This research will generate test cases for *multi-way interactions*, *i.e.*, longer  
21 paths in an EIG. For example, a 3-way test case is  $\langle e_1; e_2; e_3 \rangle$ ; a 4-way test case is  $\langle e_1; e_2; e_7;$   
22  $e_2 \rangle$ . Because EIG nodes do not represent events to open or close menus, or open windows, the  
23 sequences obtained from the EIG may not be executable. At execution time, other events needed to  
24 reach the EIG events are automatically generated, yielding an executable test case [49]. To allow a  
25 clean application exit, a test case is also automatically augmented with additional events that close  
26 all open modal windows before the test case terminates.

27 The function notation  $S_j = e_x(S_i)$  denotes that  $S_j$  is the state resulting from the execution of



1 event  $e_x$  in state  $S_i$ . If  $e_1$  and  $e_2$  are two different events in a GUI's EIG,  $(e_1, e_2)$  is an edge, and  
 2  $S_0 \in S_I$  is the initial state of the GUI, then  $e_1(S_0)$  is the GUI state after performing  $e_1$ ,  $e_2(S_0)$   
 3 is the GUI state after performing  $e_2$ , and  $e_2(e_1(S_0))$  is the GUI state after performing the *event*  
 4 *sequence*  $\langle e_1; e_2 \rangle$ .

## 5 4 Event Semantic Interaction Graph

6 The new feedback-based technique is based on our ability to identify sets of events that need to be  
 7 tested together in multi-way interactions. Because each event is executed using its corresponding  
 8 event handler, one could hypothesize that an event  $e_x$  whose event handler updates code element(s)  
 9 (e.g., variables, object fields) that are used (directly or indirectly) by another event  $e_y$ 's handler  
 10 could potentially influence its execution;  $e_x$  and  $e_y$  are good candidates to be tested together.  
 11 For example, consider the event handlers for the events  $e_1$  and  $e_2$  shown in Figure 3. As these  
 12 event handlers interact via the variable *currentTool*, the events  $e_1$  and  $e_2$  should be tested together.  
 13 Similarly, events  $e_3$  and  $e_4$  interact via *curZoom* and should be tested together. However, because  
 14 the handlers for  $e_1$  and  $e_3$  do not interact, these events need not be tested together.

<b><math>e_1</math>:: select ellipse tool</b>
public void ellipsePerformed (java.awt.event.ActionEvent evt){ ...; currentTool = toolEllipse; ... }
<b><math>e_2</math>:: drag mouse on canvas</b>
public void mouseDragged(java.awt.event.MouseEvent evt) { ...; currentTool.dragAction(newEvt, center); ... }
<b><math>e_3</math>:: set zoom factor to double</b>
public void zoom1Performed(java.awt.event.ActionEvent evt) { ...; curZoom = zoom1; ... }
<b><math>e_4</math>:: click left mouse button on canvas</b>
public void mousePressed(java.awt.event.MouseEvent evt) { ... if (currentTool == toolZoom){ // if the zoom tool is being used int temp = toolZoom.getZoom(); // current zoom level if (SwingUtilities.isLeftMouseButton(evt){ switch (temp){ case 1: zoom2.setBG(pColor); curZoom = zoom2; case 2: zoom3.setBG(pColor); curZoom = zoom3;... } } } theCanvas.repaint();... }

Figure 3: Example Event Handlers

1 One may employ a variety of static program-analysis techniques to identify such interactions  
2 [64]. They can certainly be used in this example. However, the limitations of static analysis in the  
3 presence of multi-language GUI implementations, callbacks for event handlers, virtual function  
4 calls, reflection, and multi-threading are well known [43, 64]; these limitations complicate the  
5 application of static analysis. Moreover, because most GUI applications employ a large number of  
6 library elements (*e.g.*, Java Swing), source code (an essential component for most static analysis  
7 techniques) may not be available for parts of the GUI.

8 This research avoids static analysis; instead it approximates the identification of interactions  
9 between event handlers by analyzing feedback from the run-time state of the GUI on an initial test  
10 suite. Recall that testing all two-way interactions between events is already quite practical with  
11 the smoke test suite; we treat this suite as a starting point to collect the feedback. The remaining  
12 question, addressed in this section, is: *What dynamic behavior constitutes event interaction?*

13 Consider the example shown in Figure 4. The top-left shows the *initial state* ( $S_0$ ) of an appli-  
14 cation. After an event  $e_1$  (*Select Ellipse tool*; event handler shown in Figure 3) is executed,  
15 the GUI changes its state to the one shown in the top-right ( $e_1(S_0)$ ). In this state, the “ellipse tool”  
16 remains selected. Starting from  $S_0$ , one can execute another event  $e_2$  (*Drag mouse on canvas*) and  
17 obtain the state shown in the bottom-left ( $e_2(S_0)$ ); an area of the canvas has been selected. If, how-  
18 ever, the sequence  $\langle e_1; e_2 \rangle$  is executed in  $S_0$ , a new state ( $e_2(e_1(S_0))$ ), shown in the bottom-right  
19 is obtained; an ellipse has been created. This execution is equivalent to the execution of event  $e_2$   
20 in the state  $e_1(S_0)$ . The sequence  $\langle e_1; e_2 \rangle$  produces a GUI state that shows the influence of  $e_1$   
21 on  $e_2$ . Hence,  $e_1$  and  $e_2$  are good candidates to be tested together in longer sequences to check for  
22 interaction problems. The code for  $e_1$  and  $e_2$  (previously seen in Figure 3) shows that they do in  
23 fact interact.

24 The remaining problem is to automatically compute the run-time relationship between event  
25  $e_1$  and  $e_2$  of Figure 4. We use four state descriptions for this computation:  $S_0$ ,  $e_1(S_0)$ ,  $e_2(S_0)$ , and  
26  $e_2(e_1(S_0))$ . More specifically, to compute this relationship, we need to find at least one *new* widget  
27  $w$  with property  $p$  and value  $v$  in state  $e_2(e_1(S_0))$ , *i.e.*, it is created by event sequence  $\langle e_1; e_2 \rangle$ ;

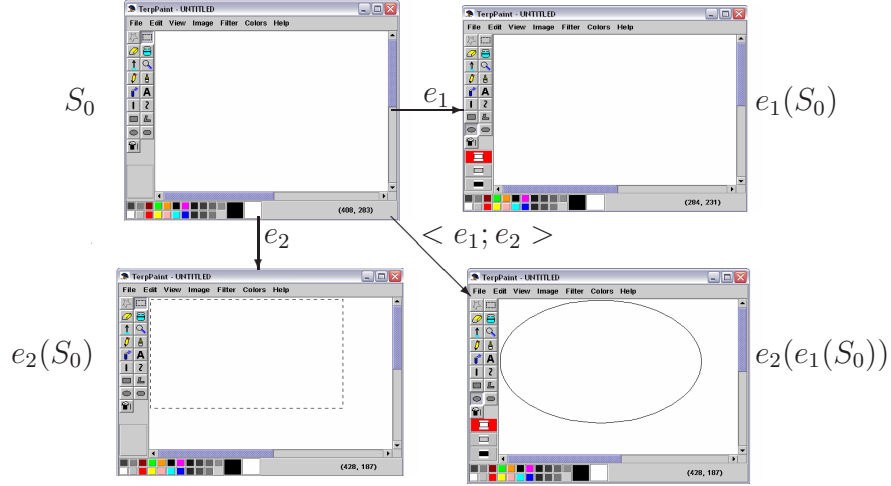


Figure 4: Execution of Events  $e_1$  (*Select Ellipse tool*) and  $e_2$  (*Drag mouse on canvas*)

1 but it does not exist in state  $S_0$  and could not be created by either  $e_1$  or  $e_2$  individually, *i.e.*, no triple  
 2 involving widget  $w$  exists in any of the states  $S_0$ ,  $e_1(S_0)$  and  $e_2(S_0)$ . This statement can be formally  
 3 described as a predicate:  $\exists w \in W, \exists p \in P_w, \exists v \in V_p, \forall \bar{p} \in P_w, \forall \bar{v} \in V_p, s.t.^4 ((w, \bar{p}, \bar{v}) \notin$   
 4  $S_0) \wedge ((w, \bar{p}, \bar{v}) \notin e_1(S_0)) \wedge ((w, \bar{p}, \bar{v}) \notin e_2(S_0)) \wedge ((w, p, v) \in e_2(e_1(S_0)))$ . If this predicate  
 5 evaluates to TRUE for a widget (in our case the ellipse), then we have discovered a case where  $e_1$   
 6 influences  $e_2$ .

7 It is quite straightforward to encode such a predicate in a high-level programming language,  
 8 such as Java. The implementation would loop through the state triples and stop when one widget  
 9 satisfying the predicate is detected.

10 The reader should note that the above predicate is necessary for the computation of the rela-  
 11 tionship between  $e_1$  and  $e_2$ . At first glance, the reader might be tempted to think that checking  
 12 state non-equivalence would be sufficient to identify interacting events, *i.e.*, by using a predi-  
 13 cate  $\mathcal{P}$  such as  $(e_1(S_0) \neq e_2(e_1(S_0))) \vee (e_2(S_0) \neq e_1(e_2(S_0)))$ . However, this is not the case.  
 14 Consider an example of two non-interacting events,  $e_x$  and  $e_y$ , which toggle the states of two inde-  
 15 pendent check-box widgets  $\square_x$  and  $\square_y$ , respectively. Starting in a state  $S_0 = \{\square_x, \square_y\}$ , *i.e.*, both  
 16 boxes unchecked, each event would “check” its corresponding check-box, *i.e.*,  $e_x(S_0) = \{\surd_x, \square_y\}$ ,  
 17  $e_y(S_0) = \{\square_x, \surd_y\}$ , and  $e_y(e_x(S_0)) = \{\surd_x, \surd_y\}$ . Even though  $\mathcal{P}$  would evaluate to TRUE for this

<sup>4</sup>Notation for “such that”

1 example, events  $e_x$  and  $e_y$  are non-interacting and need not be tested together.

2 Additional complexity arises from the nature of modal windows in GUIs. Modal windows  
3 create special situations due to the presence of termination events. For example, instead of being  
4 in the main window, had  $e_1$  and  $e_2$  in Figure 4 been associated with widgets contained in a modal  
5 window with termination event TERM, states  $e_1(S_0)$ ,  $e_2(S_0)$ , and  $e_2(e_1(S_0))$  might not contain the  
6 necessary information needed to compute the predicate. This is because user actions in modal  
7 windows do not cause immediate state changes; they typically take effect after a termination event  
8 has been executed. Hence, each of the states  $e_1(S_0)$ ,  $e_2(S_0)$ , and  $e_2(e_1(S_0))$  must be collected  
9 after the execution of the termination event TERM. Similarly, problems arise when  $e_1$  and  $e_2$  are in  
10 two *different* modal windows;  $e_1$  is in a modal window but  $e_2$  is in a modeless window;  $e_1$  is in a  
11 modal window whereas  $e_2$  is in its parent window. All these situations require special handling.

12 It turns out that the example illustrated in Figure 4 is just one *case* of how the GUI state may be  
13 used to pinpoint interactions between event handlers – there are several more cases, each requiring  
14 a different predicate. Because of the need to define precise predicates for all these cases and for  
15 special handling of modal windows, we take a two dimensional approach. We first define six<sup>5</sup>  
16 predicates in one *context*, *i.e.*, where  $e_1$  and  $e_2$  are system-interaction events in modeless windows;  
17 this situation is called *Context 1*. We will use the notation  $\mathcal{P}_{n(m)}(e_1, e_2)$  to represent a predicate  
18 for case  $n$  in context  $m$ . We then define two additional contexts; together, the six cases and three  
19 contexts yield  $6 \times 3 = 18$  situations for computing run-time relationships between events.

20 **Case 1:**  $\mathcal{P}_{1(1)}(e_1, e_2) = \exists w \in W, p \in P_w, v \in V_p, v' \in V_p, s.t. ((v \neq v') \wedge ((w, p, v) \in \{S_0 \cap$   
21  $e_1(S_0) \cap e_2(S_0)\}) \wedge ((w, p, v') \in e_2(e_1(S_0))))$ ; there is at least one widget  $w$  with property  $p$  with  
22 initial value  $v$  (hence the triple  $(w, p, v)$  is in  $S_0$ ), which is not affected by the individual events  $e_1$   
23 or  $e_2$  (the triple is also in  $e_1(S_0)$  and  $e_2(S_0)$ ); however, it is modified when the sequence  $\langle e_1; e_2 \rangle$   
24 is executed, *i.e.*, the value of  $w$ 's property  $p$  changes from  $v$  to  $v'$ .

25 Figure 5 gives an example of **Case 1**. This is a “GUI Demo” application with several widgets.  
26 The `Fill with color` checkbox fills the currently selected shape (highlighted with a deep

---

<sup>5</sup>We have chosen to present only these six cases because we encountered them numerous times in our work on GUI testing. These cases are not exhaustive and we will continue to add new cases, as and when needed, in the future.

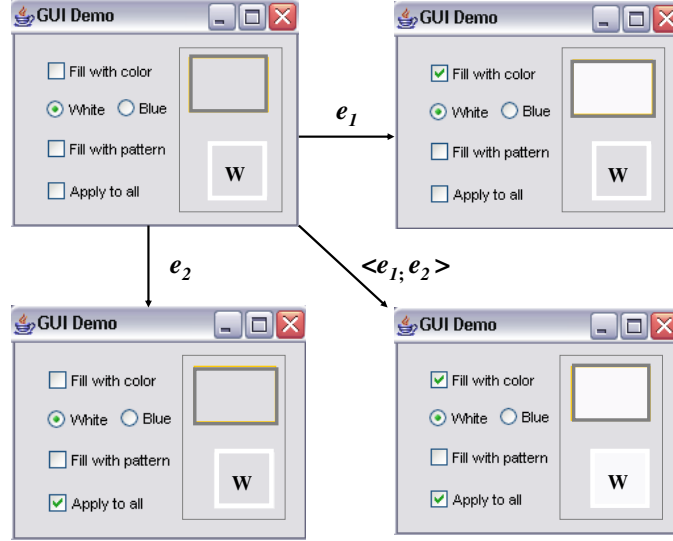


Figure 5: **Case 1**:  $e_1$ : *Check Fill with color*;  $e_2$ : *Check Apply to all*

1 grey border) with the chosen color determined by the radio buttons White and Blue. Checkbox  
 2 Fill with pattern determines whether to fill the selected shape with a pattern. Checking  
 3 Apply to all sets all shapes in the right panel with the same color and pattern.

4 For the purpose of **Case 1**,  $e_1$  is *Check Fill with color* and  $e_2$  is *Check Apply to*  
 5 *all*. The initial state has the rectangle widget selected and color is set to white. The square  
 6 widget (marked with **W**) is not modified by  $e_1$  or  $e_2$  individually; however, the event sequence  
 7  $\langle e_1; e_2 \rangle$  fills the square with the white color. Hence  $\mathcal{P}_{1(1)}(e_1, e_2)$  evaluates to TRUE – **Case 1** is  
 8 applicable here and  $e_1$  is ESI related to  $e_2$  because  $e_1$  influences  $e_2$  and their combination modifies  
 9 the previously unmodified widget **W**.

10 **Case 2**:  $\mathcal{P}_{2(1)}(e_1, e_2) = \exists w \in W, p \in P_w, v \in V_p, v' \in V_p, v'' \in V_p, s.t. ((v \neq v') \wedge (v' \neq$   
 11  $v'') \wedge ((w, p, v) \in \{S_0 \cap e_2(S_0)\}) \wedge ((w, p, v') \in e_1(S_0)) \wedge ((w, p, v'') \in e_2(e_1(S_0))))$  there is at  
 12 least one widget  $w$  with property  $p$  that has an initial value  $v$ , which is not modified by the event  
 13  $e_2$ ; it is modified by  $e_1$ ; however, it is modified differently by the sequence  $\langle e_1; e_2 \rangle$ .

14 An example of **Case 2** using the “GUI Demo” application is given in Figure 6, where  $e_1$  now  
 15 represents *Check Fill with color* and  $e_2$  is *Click radio button Blue*. The initial state has  
 16 the rectangle selected and color is set to white. Individually, in this initial state, event  $e_1$  fills

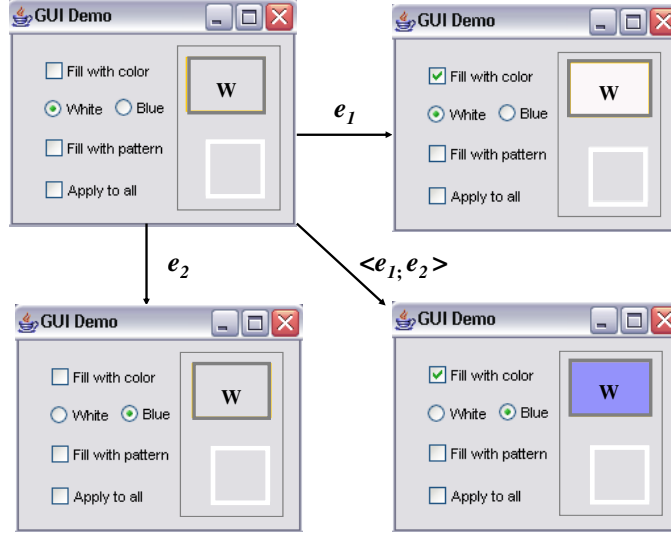


Figure 6: **Case 2:**  $e_1$ : Check Fill with color;  $e_2$ : Click radio button Blue

1 the rectangle with the white color; event  $e_2$  sets the current color to blue. However, executing

2  $\langle e_1; e_2 \rangle$  now fills the rectangle with the color blue. Hence  $\mathcal{P}_{2(1)}(e_1, e_2)$  evaluates to TRUE –

3 **Case 2** applies here as  $e_1$  influences  $e_2$  execution; the widget (marked with **W**) is modified by  $e_1$ ;

4 it is not modified by  $e_2$ ; however, it is modified differently by the sequence  $\langle e_1; e_2 \rangle$ .

5 **Case 3:**  $\mathcal{P}_{3(1)}(e_1, e_2) = \exists w \in W, p \in P_w, v \in V_p, v' \in V_p, v'' \in V_p, s.t. ((v \neq v') \wedge (v' \neq$

6  $v'') \wedge ((w, p, v) \in \{S_0 \cap e_1(S_0)\}) \wedge ((w, p, v') \in e_2(S_0)) \wedge ((w, p, v'') \in e_2(e_1(S_0)))$  there is at

7 least one widget  $w$  with property  $p$  that has an initial value  $v$ , which is not modified by the event

8  $e_1$ ; it is modified by  $e_2$ ; however, it is modified differently by the sequence  $\langle e_1; e_2 \rangle$ . Note that

9 this case is different from Case 2 because the event sequence remains the same, *i.e.*,  $e_1$  is executed

10 before  $e_2$ .

11 An example of **Case 3** using the “GUI Demo” application is given in Figure 7, where  $e_1$  now

12 represents *Click radio button Blue* and  $e_2$  is *Check Fill with color*. The initial state has

13 the rectangle selected and color is set to white. Individually, in this initial state, event  $e_1$  sets

14 the current color to blue; event  $e_2$  fills the rectangle with the white color. However, executing

15  $\langle e_1; e_2 \rangle$  now fills the rectangle with the color blue. Hence  $\mathcal{P}_{3(1)}(e_1, e_2)$  evaluates to TRUE –

16 **Case 3** applies here as  $e_1$  influences  $e_2$  execution; the widget (marked with **W**) is not modified by

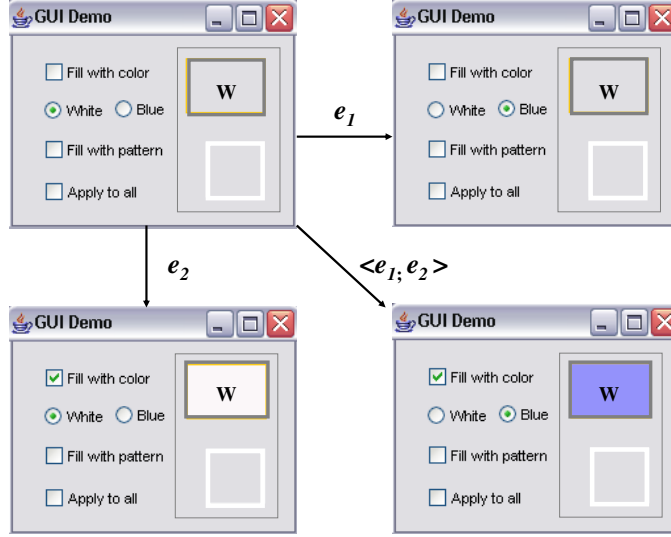


Figure 7: **Case 3:**  $e_1$ : Click radio button Blue;  $e_2$ : Check Fill with color

$e_1$ ; it is modified by  $e_2$ ; however, it is modified differently by the sequence  $\langle e_1; e_2 \rangle$ .

**Case 4:**  $\mathcal{P}_{4(1)}(e_1, e_2) = \exists w \in W, p \in P_w, v \in V_p, v' \in V_p, v'' \in V_p, \bar{v} \in V_p, s.t. ((v \neq v') \wedge (v \neq v'') \wedge (v'' \neq \bar{v}) \wedge ((w, p, v) \in S_0) \wedge ((w, p, v') \in e_1(S_0)) \wedge ((w, p, v'') \in e_2(S_0)) \wedge ((w, p, \bar{v}) \in e_2(e_1(S_0))))$ ; there is at least one widget  $w$  with property  $p$  that has an initial value  $v$ , which is modified by individual events  $e_1$  and  $e_2$ ; however, it is modified differently by the sequence  $\langle e_1; e_2 \rangle$ .

Figure 8 shows one example of this case using the “GUI Demo” application. In this example, the initial state has Fill with color checked, white is set to be the current color and the rectangle is selected. Event  $e_1$  here is Click radio button Blue and  $e_2$  is Check Fill with pattern that fills the current shape with a pattern. Events  $e_1$  and  $e_2$  modify the rectangle individually; however, executing  $\langle e_1; e_2 \rangle$  now modifies the rectangle differently. Therefore,  $e_1$  influences  $e_2$ , *i.e.*, resulting in different modification of the existing widget (marked with **W**), and **Case 4** applies because  $\mathcal{P}_{4(1)}(e_1, e_2)$  evaluates to TRUE.

The above four cases all handle widgets that persist across the four states being considered, *i.e.*,  $S_0, e_1(S_0), e_2(S_0),$  and  $e_2(e_1(S_0))$ . In many cases, event execution “creates” new widgets, *e.g.*, by opening menus; the next case handles newly created widgets.

1 **Case 5:**  $\mathcal{P}_{5(1)}(e_1, e_2) = \exists w \in W, p \in P_w, v \in V_p, v' \in V_p, s.t. ((v \neq v') \wedge ((w, p, v) \in e_x(S_0)) \wedge$   
2  $((w, p, v) \notin S_0) \wedge ((w, p, v') \in e_2(e_1(S_0))))$ ; there is at least one *new* widget  $w$  with property  $p$   
3 and value  $v$  in  $e_x(S_0)$ , *i.e.*, it was created by event  $e_x$  (either  $e_1$  or  $e_2$ ) but did not exist in state  $S_0$ ;  
4 it was created by the sequence  $\langle e_1; e_2 \rangle$  but with a different value for some property.

5 Figure 9 shows an example for **Case 5** using the “Radio Button Demo” application in-  
6 troduced earlier. The initial state has the Circle and None radio button selected and an empty  
7 Rendered Shape panel. In this initial state, event  $e_2$  selectes Square; event  $e_6$  clicks the  
8 button Create Shape and creates a circle in the Rendered Shape panel. However, when  
9 executing  $\langle e_2; e_6 \rangle$  in the initial state, a square is created in the Rendered Shape panel.  
10 Therefore,  $e_2$  influences  $e_6$ ;  $\mathcal{P}_{5(1)}(e_2, e_6)$  evaluates to TRUE – **Case 5** is applicable here.

11 It turns out that in this simple application, **Case 5** also applies to  $\langle e_6; e_2 \rangle$ ,  $\langle e_3; e_6 \rangle$  and  
12  $\langle e_6; e_3 \rangle$ .

13 A common occurrence of event interaction in GUIs is enabling/disabling widgets, which may  
14 be modeled as the widget’s ENABLED property being set to TRUE or FALSE.

15 **Case 6:**  $\mathcal{P}_{6(1)}(e_1, e_2) = \exists w \in W, ENABLED \in P_w, TRUE \in V_{ENABLED}, FALSE \in V_{ENABLED}, s.t. (((w, ENABLED, FALSE)$   
16  $S_0) \wedge ((w, ENABLED, TRUE) \in e_1(S_0)) \wedge EXEC(e_2, w))$ ; there exists at least one widget  $w$  that was dis-  
17 abled in  $S_0$  but enabled by  $e_1$ . Event  $e_2$  is performed on  $w$ , represented by a predicate  $EXEC(e_2, w)$ .

18 Using the Radio Button Demo application, we see in Figure 10 that  $e_3$  enables  $e_5$ . Hence  
19  $\mathcal{P}_{6(1)}(e_3, e_5)$  evaluates to TRUE – **Case 6** applies.

20 As mentioned earlier, the second dimension of our definitions are contexts. This is because  
21 modal windows create special situations for Cases 1 through 6 due to the presence of termination  
22 events. User actions in these windows do not cause immediate state changes; they typically take  
23 effect after a termination event has been executed, leading to *contexts 2 and 3*.

24 **Context 2:** If both  $e_1$  and  $e_2$  are associated with widgets that are contained in one modal window  
25 with termination event TERM, then the definitions of  $e_1(S_0)$ ,  $e_2(S_0)$ , and  $e_2(e_1(S_0))$  are modified  
26 as follows:  $e_1(S_0)$  is the state of the GUI after the execution of the event sequence  $\langle e_1; TERM \rangle$ ,  
27  $e_2(S_0)$  is the state of the GUI after the execution of the event sequence  $\langle e_2; TERM \rangle$ , and



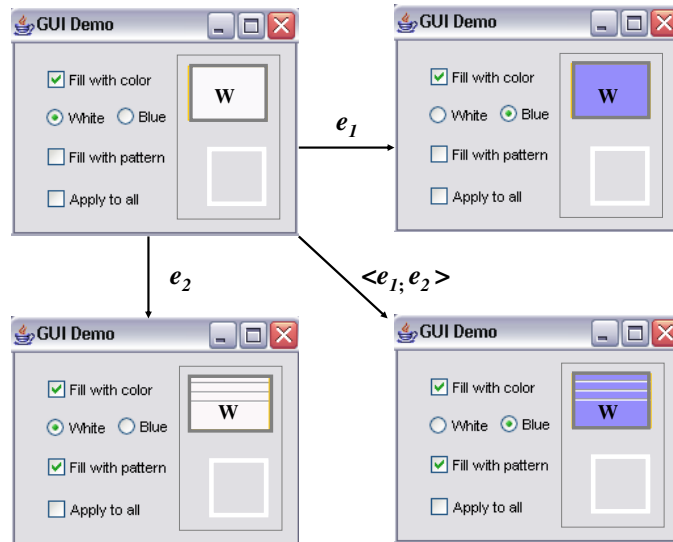


Figure 8: **Case 4:**  $e_1$ : Click radio button Blue;  $e_2$ : Check Fill with pattern

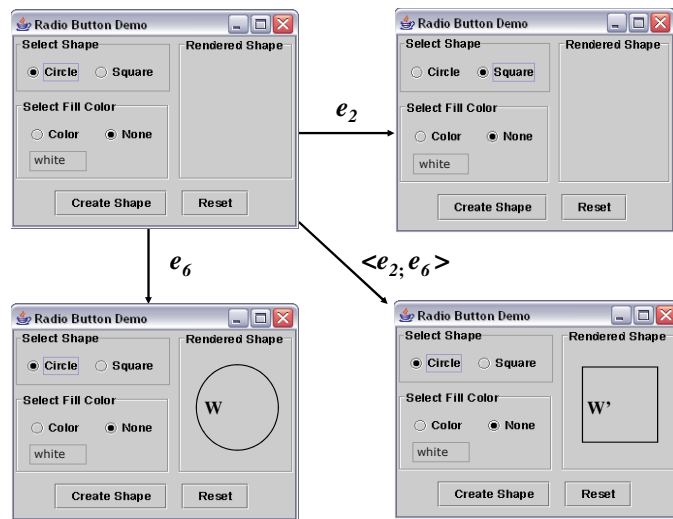


Figure 9: **Case 5:**  $e_2$ : Select Square;  $e_6$ : Click button Create Shape

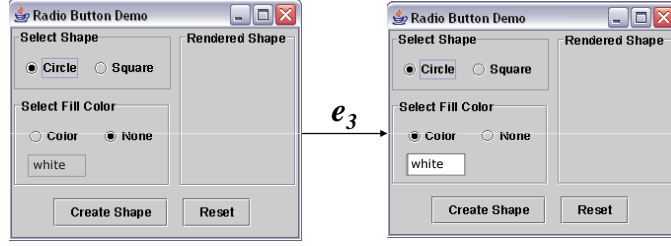


Figure 10: **Case 6:**  $e_3$ : *Select Color*;  $e_5$ : *Input color*

1  $e_2(e_1(S_0))$  is the state of the GUI after the execution of the event sequence  $\langle e_1; e_2; \text{TERM} \rangle$ .  
 2 All the predicates defined in Cases 1 through 6 apply, using these modified definitions, for  $e_1$  and  
 3  $e_2$  in the same modal window. The notation used for these predicates when applied in Context 2 is  
 4  $\mathcal{P}_{n(2)}(e_1, e_2)$ , where  $n$  is the case number.

5 **Context 3:** If  $e_1$  is associated with a widget contained in a modal window with termination event  
 6  $\text{TERM}$ , and  $e_2$  is associated with a widget contained in the modal window's *parent* window (*i.e.*,  
 7 the window that was used to open the modal window) then  $e_1(S_0)$  is the state of the GUI after the  
 8 execution of the event sequence  $\langle e_1; \text{TERM} \rangle$ ,  $e_2(S_0)$  is the state of the GUI after the execution  
 9 of the event  $e_2$ , and  $e_2(e_1(S_0))$  is the state of the GUI after the execution of the event sequence  
 10  $\langle e_1; \text{TERM}; e_2 \rangle$ . All the predicates defined in Cases 1 through 6 apply. The notation used for  
 11 these predicates when applied in Context 3 is  $\mathcal{P}_{n(3)}(e_1, e_2)$ , where  $n$  is the case number.

12 We are now ready to formally define the ESI relationship. There is an *Event Semantic Interac-*  
 13 *tion* relationship between two events  $e_1$  and  $e_2$  iff  $\mathcal{P}_{1(1)}(e_1, e_2) \vee \mathcal{P}_{2(1)}(e_1, e_2) \vee \dots \vee \mathcal{P}_{6(1)}(e_1, e_2) \vee$   
 14  $\mathcal{P}_{1(2)}(e_1, e_2) \vee \mathcal{P}_{2(2)}(e_1, e_2) \vee \dots \vee \mathcal{P}_{6(2)}(e_1, e_2) \vee \mathcal{P}_{1(3)}(e_1, e_2) \vee \mathcal{P}_{2(3)}(e_1, e_2) \vee \dots \vee \mathcal{P}_{6(3)}(e_1, e_2)$ . That  
 15 is, at least one of the predicates in Cases 1 through 6 evaluates to  $\text{TRUE}$  in at least one context;  
 16 this relationship is written as  $e_1 \xrightarrow{n(m)} e_2$ , where the number  $n$  is one of the case numbers 1 through  
 17 6;  $m$  is the context number. If multiple cases apply, then one of the case numbers is used. Due to  
 18 the specific ordering of the events in the sequence  $\langle e_1; e_2 \rangle$ , the ESI relationship is not symmet-  
 19 ric. As demonstrated earlier, for our Radio Button Demo application,  $e_2 \xrightarrow{5(1)} e_6$ ,  $e_6 \xrightarrow{5(1)} e_2$ ,  
 20  $e_3 \xrightarrow{5(1)} e_6$ ,  $e_6 \xrightarrow{5(1)} e_3$ , and  $e_3 \xrightarrow{6(1)} e_5$ .

21 Once all of the cases have been implemented, the feedback-based process execution is straight-

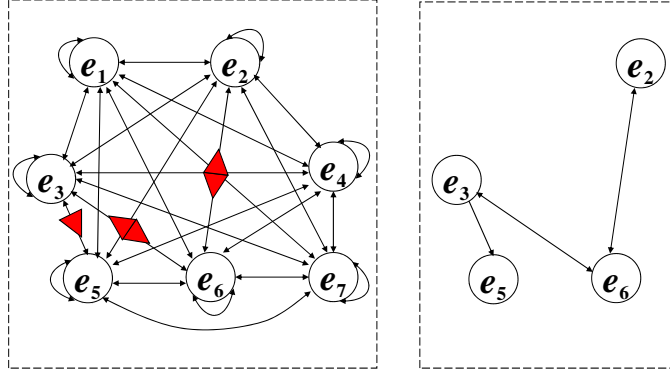


Figure 11: EIG with marked ESI relationships and ESIG for “Radio Button Demo” GUI

1 forward. The steps of the execution are as follows.

- 2 1. The seed suite consisting of all 2-way interactions  $\langle e_x; e_y \rangle$  between GUI events is exe-  
 3 cuted on the software in state  $S_0$ ; these test cases are simple enumerations of all EIG edges.  
 4 All events  $e_y$  are also executed in  $S_0$ . The state information  $e_x(S_0)$ ,  $e_y(S_0)$ ,  $e_y(e_x(S_0))$  is  
 5 collected and stored.
- 6 2. The above predicates are evaluated for each pair of system-interaction events in the EIG that  
 7 are either (1) directly connected by an edge (Context 1) or (2) connected by a path that does  
 8 not contain any intermediate system-interaction events (contexts 2 and 3), *i.e.*, there is at  
 9 least one termination event that closes a modal window on this path. If one of the predicates  
 10 evaluates to TRUE, the two events are ESI-related.

11 Once all the ESIs in a GUI have been identified, a graph model called the ESI graph (ESIG) is  
 12 created. The ESIG contains nodes that represent events; a directed edge from node  $n_x$  to  $n_y$  shows  
 13 that there is an ESI relationship from the event represented by  $n_x$  to the event represented by  $n_y$ .  
 14 The EIG annotated with the five ESI relationships found in Radio Button Demo are shown in  
 15 Figure 11; the ESIG (shown on the right) is a subgraph of the EIG.

16 The ESIG may be traversed using a modified version of the GenTestCases algorithm dis-  
 17 cussed in Section 3. The differences are that (1) an ESIG may contain multiple connected compo-  
 18 nents in which case the event sequences are generated for each component separately, and (2) the  
 19 length of the obtained sequences is now a tunable parameter instead of a fixed number 2. Study 1

1 in the next section uses values 3, 4, and 5 for this parameter.

2 Our new implementation of the `GenTestCases` algorithm is based on the adjacency matrix  
3 representation of directed graphs. The key idea used in the implementation is that if we start with  
4 a 0/1 adjacency matrix representation of the ESIG, and take that matrix to the  $(N - 1)^{th}$  power,  
5 the  $(i, j)$  entry in the resulting matrix is the number of paths of length  $N$  from node  $i$  to node  $j$   
6 (recall that the length is measured in number of nodes encountered along the path). In the trivial  
7 case,  $N = 2$  will return the input matrix – the  $(i, j)$  entry is either 0 or 1, *i.e.*, the number of length  
8 2 paths from node  $i$  to node  $j$ . For  $N = 3$ , the  $(i, j)$  entry in the result matrix is the number of all  
9 length 3 paths from node  $i$  to  $j$ .

10 Because we want to output actual test cases, not just count them, we use a variation of the above  
11 approach. The only difference is that instead of just counting the paths, our implementation keeps  
12 track of all the actual paths themselves. For this we had to modify the matrix multiplication algo-  
13 rithm and the adjacency matrix representation. The adjacency matrix is modified so that instead of  
14 0/1, the  $(i, j)$  entry of the matrix is a list of paths from  $i$  to  $j$ . The matrix multiplication algorithm  
15 is modified so that instead of multiplying and adding entries, we instead concatenate pairs of paths  
16 together and union all of them (respectively) to eliminate duplicates. The final matrix entries are  
17 paths of specific lengths, *i.e.*, test cases of specific lengths.

18 Much of the functionality needed for this test-case generation approach is implemented in the  
19 *Mathematica* package. The function `MatrixPower[mat, n]`, returns the  $n^{th}$  matrix power  
20 of matrix `mat`. If the matrix is encoded as a 0/1 adjacency list, `MatrixPower` returns a matrix  
21 in which each entry  $(i, j)$  is a count of the number of test cases of a specific length from node  
22  $i$  to node  $j$ . We used simple rewriting rules (provided as a `Replace` function denoted by the  
23 operator “/.”) built into *Mathematica* to alter `MatrixPower` – we replaced the `Times` (*i.e.*,  
24 multiply) operator with `Join` (*i.e.*, list concatenation) and the `Plus` operator with `Union`. We  
25 thus generate all length  $N - 1$  test cases by using the *Mathematica* command:

```
26 MatrixPower[matrix, N]/.{Power[x_, 2]->Join[x, x], Plus->Union,  
27 Times->Join}
```

1 with an appropriate representation of the matrix.

2 Because our ESIG contains very few edges and nodes, the adjacency matrix tends to be sparse,  
3 leading to fast computations. For increased efficiency, *Mathematica*'s `MatrixPower` uses the  
4 built-in `SparseArray` structure.

5 For our example ESIG of Figure 11, the test cases of length 3 are  $\langle e_2; e_6; e_2 \rangle$ ,  $\langle e_2; e_6; e_3 \rangle$   
6 ,  $\langle e_3; e_6; e_3 \rangle$ ,  $\langle e_3; e_6; e_2 \rangle$ ,  $\langle e_6; e_3; e_6 \rangle$ , and  $\langle e_6; e_2; e_6 \rangle$ . Note that because there is no  
7 non-determinism in the test-case generation algorithm, there is a unique ESIG-based test suite of  
8 specific-length test cases for an application.

## 9 **5 Study 1: Evaluating the Feedback-based Technique on Fielded** 10 **Applications**

11 The test cases obtained from the modified `GenTestCases` algorithm can be generated and ex-  
12 ecuted automatically on the GUI. The only unavailable part is the *test oracle*, a mechanism that  
13 determines whether an AUT executed correctly for a test case. In this first study, an AUT is consid-  
14 ered to have *passed* a test case if it did not “crash” (terminate unexpectedly or throw an uncaught  
15 exception) during the test case’s execution; otherwise it *failed*. Such crashes may be detected au-  
16 tomatically by the script used to execute the test cases. The EIG and ESIG, and their respective  
17 test cases may also be obtained automatically. Hence, the entire end-to-end feedback-based GUI  
18 testing process for “crash testing” could be executed without human intervention. Note that, in the  
19 next section (Study 2), this work is extended by employing a more “powerful” test oracle to detect  
20 additional failures.

21 Implementation of the crash testing process included setting up a database for text-field values.  
22 Since the overall process needed to be fully automatic, a database containing one instance for each  
23 of the text types in the set  $\{\textit{negative number}, \textit{real number}, \textit{long file name}, \textit{empty string}, \textit{special}$   
24  $\textit{characters}, \textit{zero}, \textit{existing file name}, \textit{non-existent file name}\}$  was used. Note that if a text field is  
25 encountered in the GUI, one instance for each text type is tried in succession.

1 This process provided a starting point for a feasibility study to evaluate the ESIG-generated test  
2 cases. The following questions needed to be answered to determine the usefulness of the overall  
3 feedback-based process:

4 **S1Q1:** How many test cases are required to test two-way interactions in an EIG? How does this  
5 number grow for 3-, 4-, ..., n-way interactions?

6 **S1Q2:** In how many ESI relationships does a given event participate? How many test cases are  
7 required to test two-way interactions in an ESIG? How does this number grow for 3-, 4-, ..., n-way  
8 interactions?

9 **S1Q3:** How do the ESIG- and EIG-generated test suites compare in terms of fault-detection effec-  
10 tiveness? Do the former detect faults that were not detected by the latter?

11 More specifically, the following process was used for this study:

- 12 1. Select software subjects with GUI front-ends.
- 13 2. Generate a seed test suite using EIGs.
- 14 3. Execute the seed test suite. Report crashes.
- 15 4. During execution, construct the ESIG. Use the ESIG to generate additional test cases.
- 16 5. Execute the new test cases. Report crashes.

17 To answer the above questions while minimizing threats to external validity, this study was con-  
18 ducted using four extremely popular GUI-based open-source software (OSS) applications down-  
19 loaded from SourceForge. The fully-automatic crash testing process was executed on them and  
20 the cause (*i.e.*, the *fault*) of each crash in the source code was determined.

21 **STEP 1: Selection of subject applications.** Four popular GUI-based OSS (CrosswordSage 0.3.5,  
22 FreeMind 0.8.0, GanttProject 2.0.1, JMSN 0.9.9b2) were downloaded from SourceForge. These  
23 applications have been used in our previous experiments [80, 86]; details of why they were cho-  
24 sen have been presented therein. In summary, all the applications have an active community of  
25 developers and a high all-time-activity percentile on SourceForge. Due to their popularity, these  
26 applications have undergone quality assurance before release. To further eliminate “obvious” bugs,  
27 a static analysis tool called *FindBugs* [29] was executed on all the applications; after the study, we

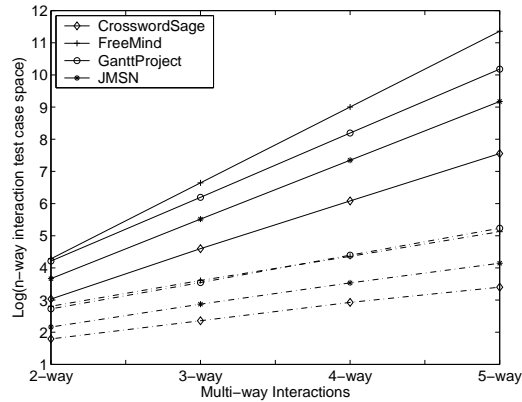


Figure 12: Test Case Space Growth

1 verified that none of our reported bugs were detected by FindBugs.

2 **STEP 2: Generation of EIGs & seed test suites.** The EIGs of all subject applications were  
 3 obtained using reverse engineering. To address **S1Q1** above, the number of test cases required  
 4 to test 2-, 3-, 4-, and 5-way interactions was computed. The result for each application is shown  
 5 as a solid line in Figure 12 (the y-axis in all these plots is a logarithmic scale). The plot shows  
 6 that the number of test cases grows exponentially with the number of interactions. The number  
 7 quickly becomes unmanageable for more than 2- and 3-way interactions. In this study, only two-  
 8 way interactions were tested by the seed test suites. The seed test suites contained 920; 51,316;  
 9 29,033; and 4634 test cases for CrosswordSage, FreeMind, GanttProject, and JMSN, respectively.

10 **STEP 3: Execution of the seed test suite.** The entire seed suite executed without any human  
 11 intervention. It executed in 0.39, 30.83, 22.89, and 2.68 hours on CrosswordSage, FreeMind,  
 12 GanttProject, and JMSN, respectively. In all, 163, 66, 14, and 34 test cases caused crashes; these  
 13 crashes were caused by 5, 4, 3, and 3 *faults* (as defined earlier) for CrosswordSage, FreeMind,  
 14 GanttProject, and JMSN, respectively. The GUI's run-time state was recorded during test execu-  
 15 tion. All faults were fixed in the applications.

16 Note that debugging and fault-fixing was necessary due to two reasons. First, had we not done  
 17 so, the longer test cases that we will generate in the next few steps may contain these short test  
 18 cases as subsequences; the longer tests may hence also crash due to the faults previously detected  
 19 by the seed suite, yielding no new useful results. Second, this is what would happen in a real

1 situation; a fault will be fixed after it was detected. However, this is a threat to internal validity  
2 because an obvious fix in one place may lead to a new fault at another place in the application. To  
3 minimize the threat, we reran the seed test suite to ensure the quality of the fixes. Of course, this  
4 does not preclude the possibility of introducing faults that are exposed by longer event sequences.  
5 To completely eliminate this threat, we later verified that the faults detected by our longer ESIG  
6 suites were not caused by these fixes.

7 **STEP 4: Generation of the ESIG.** The above feedback was used to obtain the ESIs for each  
8 application. To address **S1Q2**, the number of ESI relationships in which each event participates is  
9 shown in Figure 13. Each event in the GUI has been assigned a unique integer ID; all event IDs  
10 are shown on the x-axis. The y-axis shows the number of ESI relationships in which the event  
11 participates.

12 The result shows that certain events dominate (around 25%) the ESI relationship in GUIs.  
13 Manual examination of these “dominant” events revealed that the nature of the subject applications,  
14 *i.e.*, most of them have a single dominant object (crossword puzzle, mind map, project schedule,  
15 messenger window) that are the focus of most events, is such that several key events influence a  
16 large number of other events. In the future, we will create a classification of these dominant events.  
17 Moreover, several events participate in very few or no ESI relations. These events include parts of  
18 the `Help` menu that has no interaction with other application events, and windowing events such  
19 as scrolling for which no developer-written code exists.

20 The ESIs were used to obtain the ESIGs and, subsequently, additional test cases. The number  
21 of test cases required to test 2-, 3-, 4-, and 5-way interactions using an ESIG is shown, for each  
22 application, as a dotted line in Figure 12. This result shows that the growth of the ESIG-generated  
23 test cases appears manageable for 3-, 4-, and (given sufficient resources) 5-way interactions. They  
24 are in fact reduced from the EIG by 99.78%, 99.97%, and 99.99% for 3-, 4-, and 5-way inter-  
25 actions, respectively. In this study, test cases for 3-, 4-, and 5-way interactions were generated.  
26 The total number of test cases for these interactions was 3592, 160,629, 199,127, and 18,144 for  
27 CrosswordSage, FreeMind, GanttProject, and JMSN, respectively.



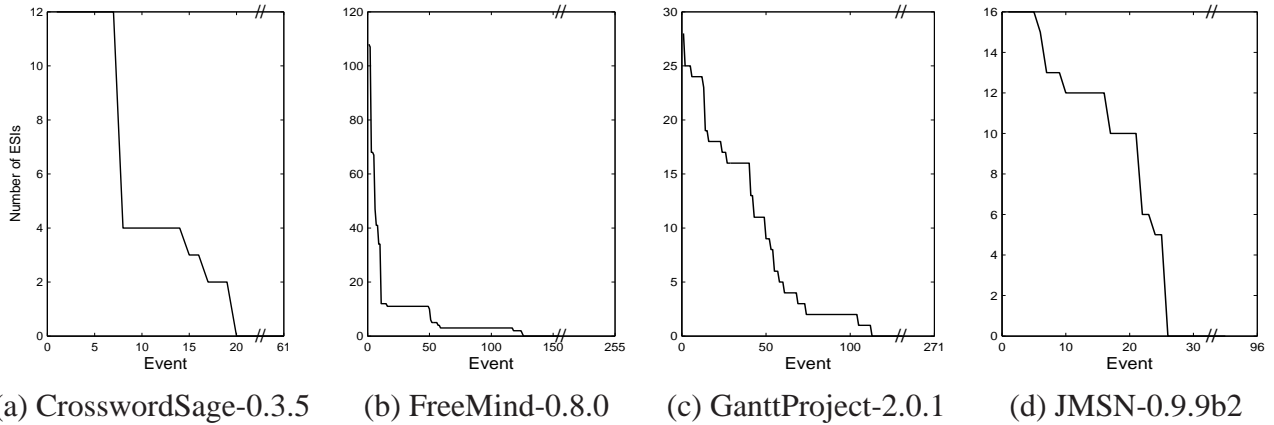


Figure 13: ESI Distribution in OSS

1 **STEP 5: Execution of the test cases.** To address **S1Q3**, all the newly-generated test cases were  
 2 executed. The execution lasted for several days. In all, 68, 157, 109, and 20 test cases caused  
 3 crashes; they were caused by 3, 3, 3, and 1 faults for CrosswordSage, FreeMind, GanttProject, and  
 4 JMSN, respectively. These faults had not been detected by the two-way test cases. We manually  
 5 verified that the faults were not introduced by our bug fixes of STEP 4. The result shows that the  
 6 ESIG-based test cases help to detect additional faults.

## 7 **6 Study 2: Digging Deeper via Seeded Faults and In-House Ap-** 8 **plications**

9 Although the previous study demonstrated the usefulness of the ESIG-based technique, it also  
 10 raised some important questions. One fundamental question that comes to mind pertains to the  
 11 cause(s) of the added effectiveness, *i.e.*, *Is the added effectiveness an incidental side-effect of the*  
 12 *events, event interactions, and lines-of-code that the ESI test cases cover and their length; or is it*  
 13 *really due to targeted testing of the identified ESI relationships?* The empirical study presented in  
 14 this section is designed specifically to address the question of how the fault-detection effectiveness  
 15 of the suite obtained by the feedback-based technique compare to that of other “similar” suites,  
 16 where similarity is quantified in terms of statement coverage, event coverage, edge coverage, and

1 size (number of test cases).

2 This question will be answered by selecting four pre-tested GUI-based applications, and gener-  
3 ating and executing 2-way EIG-based and 3-way ESIG-based test suites on them. We will generate  
4 additional test suites that are similar to the ESIG-based suite in terms of the aforementioned char-  
5 acteristics and are at least 3-way interacting, and compare their fault-detection effectiveness. Fault  
6 detection effectiveness will be measured on a per-test-suite basis in terms of number of faults de-  
7 tected. We will also study the faults, pinpointing reasons for why some of them remain undetected  
8 by our technique.

## 9 **6.1 Preparing the Subject Applications & Test Oracles**

10 Four open-source applications, called the TerpOffice suite, consisting of Paint, Present, Spread-  
11 Sheet and Word, have been selected for the study.<sup>6</sup> Table 1 shows key metrics for TerpOffice.  
12 These applications are selected very carefully for a number of reasons. In particular, to minimize  
13 threats to external validity, the selected applications are non-trivial, consisting of several GUI win-  
14 dows and widgets. For reasons described later, artificial faults were seeded in the applications –  
15 this required access to source code, bug reports, and a CVS development history. To avoid (the  
16 often difficult) distinction between GUI code and underlying “business logic,” GUI-intensive ap-  
17 plications were selected, *i.e.*, most of the source-code implemented the GUI. Finally, the tools  
18 implemented for this research, in particular for reverse engineering, are well-tuned for the Java  
19 Swing widget library – the applications had to be implemented in Java with a GUI front-end based  
20 on Swing components. As is the case with all empirical studies, the choice of subject applications  
21 introduces some significant threats to external validity of the results; these (and other) threats have  
22 been noted in Section 7.

23 For the purpose of this study, a GUI fault is a mismatch, detected by a test oracle, between  
24 the “ideal” (or expected) and actual GUI states. Hence, to detect faults, a description of ideal GUI

---

<sup>6</sup>Detailed specifications, requirements documents, source code CVS history, bug reports, and developers’ names are available at <http://www.cs.umd.edu/users/atif/TerpOffice/>.

Subjects	Windows	Widgets	LOC	Classes	Methods
Paint	16	301	11,803	330	1,253
Present	11	322	10,847	292	2,057
SpreadSheet	9	176	5,381	135	746
Word	26	617	9,917	197	1,380
TOTAL	62	1,416	38,398	954	5,436

Table 1: TerpOffice Applications

1 execution state is needed. This description is used by test oracles to detect faults in the subject  
2 applications. There are several ways to create this description. First is to manually create a formal  
3 GUI specification and use it to automatically create test oracles. Second is to use a capture/replay  
4 tool to manually develop assertions corresponding to test oracles and use the assertions as oracles  
5 to test other versions of the subject applications. Third is to develop the test oracle from a “golden”  
6 version of the subject application and use the oracle to test fault-seeded versions of the application.  
7 The first two approaches are extremely labor intensive since they require the development of a for-  
8 mal specification and the use of manual capture/replay tools; the third approach can be performed  
9 automatically and has been used in this study.

10 Several faults were seeded in each application. In order to avoid fault interaction and to sim-  
11 plify the mapping of application failure to underlying fault, multiple versions of each application  
12 were created; each version was seeded with exactly one fault. Hence, a test case detects a fault  $i$  if  
13 there is a mismatch between version  $i$  (*i.e.*, the version that was created by seeding fault  $i$ ) and the  
14 original. A mismatch is detected by comparing, between the golden and fault seeded version, the  
15 values of all the properties of all the GUI widgets being displayed, after each event.

16 The process used for fault seeding was similar to the one used in earlier work [51, 82]. Details  
17 will not be replicated here. In summary, during fault seeding, 12 classes of known faults were iden-  
18 tified, and several instances of each fault class were artificially introduced into the subject program  
19 code in source code statements that were covered by the smoke test cases, thereby ensuring that  
20 these statements were part of executable code. Care was taken so that the artificially seeded faults  
21 were similar to faults that naturally occur in real programs due to mistakes made by developers; the  
22 faults were seeded “fairly,” *i.e.*, an adequate number of instances of each fault type were seeded.

1 Several graduate students were employed to seed faults in each subject application; they created  
2 263, 265, 234, and 244 faulty versions for Paint, Present, SpreadSheet, and Word, respectively.

### 3 **6.2 Generating and Executing the ESIG-Based Test Suite**

4 The reverse engineering process was used to obtain the EIGs for the original versions of each  
5 application. The sizes of the EIGs, in terms of nodes and edges, are shown in Table 2. These  
6 numbers are important as they determine the number of generated test cases and their growth in  
7 number as test-case length increases.

	Paint	Present	SpreadSheet	Word
# <i>EIG Nodes</i>	300	321	175	616
# <i>EIG Edges</i>	21,391	32,299	6,782	28,538
# <i>ESIG Nodes</i>	102	50	45	75
# <i>ESIG Edges</i>	233	233	197	204

Table 2: ESIG vs. EIG Sizes

8 The EIGs were then used to generate all possible 2-way test cases, *i.e.*, the smoke tests. The  
9 numbers generated were exactly equal to the number of edges in the EIGs – it was quite feasible  
10 to execute such numbers of test cases in little more than a day on our 50 machines in parallel. The  
11 test cases were executed on their corresponding “correct” applications; the GUI state was collected  
12 and stored. The reader should note that it is quite impractical to generate all possible length 3 test  
13 cases for these EIGs.

14 While new software versions were being obtained (via fault seeding as discussed in Sec-  
15 tion 6.1), the 2-way EIG-based test suites and GUI state were used to obtain all possible 3-way  
16 ESIG covering test cases. The sizes of the ESIGs are shown in Table 2. The table shows that  
17 the ESIGs are much smaller than the corresponding EIGs. Due to the small number of nodes and  
18 edges, the number of 3-way covering test cases was 2531, 2080, 2069, and 2345 for Paint, Present,  
19 SpreadSheet, and Word, respectively. As noted earlier, there is a unique set of length 3 test cases  
20 for an ESIG; hence, there is a single ESIG test suite per application.

21 The 2-way EIG- and 3-way ESIG-based test cases were then executed on the fault-seeded

	Paint	Present	SpreadSheet	Word
Total Faults	263	265	234	244
2-way EIG-detected Faults	147	139	139	183
3-way ESIG-detected Faults (only new faults)	47	52	39	36

Table 3: ESIG vs. EIG Fault Detection

1 versions of the applications. The number of faults detected is shown in Table 3. Note that the last  
2 row reports the number of “new” faults detected by ESIG. This table shows that ESIG-based suites  
3 are able to detect a large number of faults missed by the EIG.

### 4 **6.3 Developing “Similar” Suites**

5 As mentioned earlier, this study required the development of several new test suites. To minimize  
6 threats to validity, the suites needed to satisfy a number of requirements, discussed next.

7 From previous studies, we know that statement, event, and EIG-edge coverage, and size (num-  
8 ber of test cases) play an important role in the fault-detection effectiveness of a test suite [81].  
9 For example, a small test suite that covers few lines of code will most likely detect fewer faults  
10 than another larger suite that covers many more lines. To allow fair comparison of fault-detection  
11 effectiveness, we needed test suites that have the *same statement, event, and edge coverage, and*  
12 *size (number of test cases)* as that of ESIG-based test suites.

13 Previous studies have also shown that long test cases (number of EIG events) fare better than  
14 short ones in terms of the number of faults that they detect [86]. Because we did not want the new  
15 suites to have any disadvantage, we ensured that all their test cases had at least 3 EIG events (note  
16 that all our ESIG test cases have exactly 3 ESIG/EIG events).

17 It is non-trivial to generate these test suites. For example, consider the problem of generating  
18 a GUI test suite that covers specific lines of code. Because of the different levels of abstraction  
19 between GUI events and code, one would need to manually examine the source code, the rela-  
20 tionship between events and underlying code, and carefully tailor each event in every test case to  
21 ensure that it covers a specific line. Because there are no automated techniques to do this task, the  
22 process will be very resource intensive.

1       Moreover, because the above criteria (same statement, event, and edge coverage, and size) may  
2 be met by a large number of test suites (with varying fault-detection effectiveness), the process of  
3 generating different suites and comparing them to the ESIG-based suites needed to be repeated  
4 several times. In this study, we generated 700 test suites per application and compared their fault-  
5 detection effectiveness to the ESIG suite.

6       GUI test cases are expensive to execute – each test case can take up to 2 minutes to execute (on  
7 average, each requires 30 seconds). Our 700 suites each for Paint, Present, SpreadSheet, and Word,  
8 contained 1,054,064; 860,324; 850,808; and 974,235 test cases, respectively; in all 3,739,431 test  
9 cases. Each test case needed to be run on each fault-seeded version; this task would have taken  
10 several years on our 50-machine cluster – an impractical task. Other researchers, who have also  
11 encountered similar issues of practicality, have circumvented this problem by creating a *test pool*  
12 consisting of a large number of test cases that can be executed in a reasonable amount of time [15].  
13 Each test case in the pool is executed only once and it’s execution attributes *e.g.*, time to execute  
14 and faults detected are recorded. Multiple test suites are created by carefully selecting test cases  
15 from this pool. Their execution is “simulated” by combining the attributes of constituent test cases  
16 using appropriate functions (*e.g.*, *set union* for faults detected). This research will also employ  
17 the test pool approach to create a large number of test suites. The test-pool-based approach will  
18 introduce some threats to validity, which we will note in Section 7.

19       Finally, we did not want to introduce any human bias when generating these test cases. We  
20 used a randomized guided mechanical process. A related approach was employed by Rothermel *et*  
21 *al.* [63] to create sequences of commands to test command-based software. In their approach, each  
22 command was executed in isolation and test cases were “assembled” by concatenating commands  
23 together in different permutations. Since GUI events (commands) enable/disable each other, most  
24 arbitrary permutations result in unexecutable sequences. Hence, we used the EIG model to obtain  
25 only executable sequences.

26       We generated test cases in batches of increasing lengths, measured in terms of the number of  
27 EIG events. We required that each EIG edge be covered by at least  $N$  test cases of a particular

1 batch. Moreover, we required that each fault-seeded statement be covered by at least  $M$  test  
2 cases of the overall pool. The test-case generation process started by generating (using a process  
3 described in the next paragraph) the batch of length-3 test cases until each EIG edge was covered  
4 by at least  $N$  test cases; they were all executed and their statement coverage was evaluated; the  
5 next (and all subsequent) batch was generated ONLY IF each fault-seeded statement was not yet  
6 covered by at least  $M$  test cases.

7 The process of generating each batch of length  $i$  test cases uses the following algorithm:

- 8 1. Initialize a `frequency` variable for each EIG edge to `zero`.
- 9 2. For each event  $e_x$  in the EIG, do
  - 10 (a) Add the single event  $e_x$  to a new empty test case  $t$ .
  - 11 (b) Form a list of all outgoing edges from  $e_x$ .
  - 12 (c) Select the edge  $(e_x, e_y)$  that has the lowest frequency, breaking ties via random selec-  
13 tion. Add  $e_y$  to the test case  $t$ .
  - 14 (d) Follow the selected edge to its destination event  $e_y$ .
  - 15 (e) Starting at  $e_y$ , recurse the `frequency`-based selection and follow-the-edge process  
16 (described in Steps 2b through 2d and this recursive step) until the desired length is  
17 obtained, adding events into the test case  $t$ .
- 18 3. Add the test case  $t$  to the suite.
- 19 4. If all EIG events have been covered and all `frequency`  $\geq N$ , stop; otherwise go to the next  
20 EIG event (via the iteration of Step 2 above).

21 The above algorithm was guaranteed to stop because all faults had been seeded in lines that  
22 were executable by the smoke tests; the count for each statement would ultimately reach  $M$  and  
23 stop. Finally, all the ESIG-based test cases were added to the pool.

24 In this study, we set  $N = 10$  and  $M = 15$ . This choice was dictated by the availability  
25 of resources. As described earlier, all the test cases needed to be executed on the fault-seeded  
26 versions of their respective application. Even with 50 machines running the test cases in parallel,

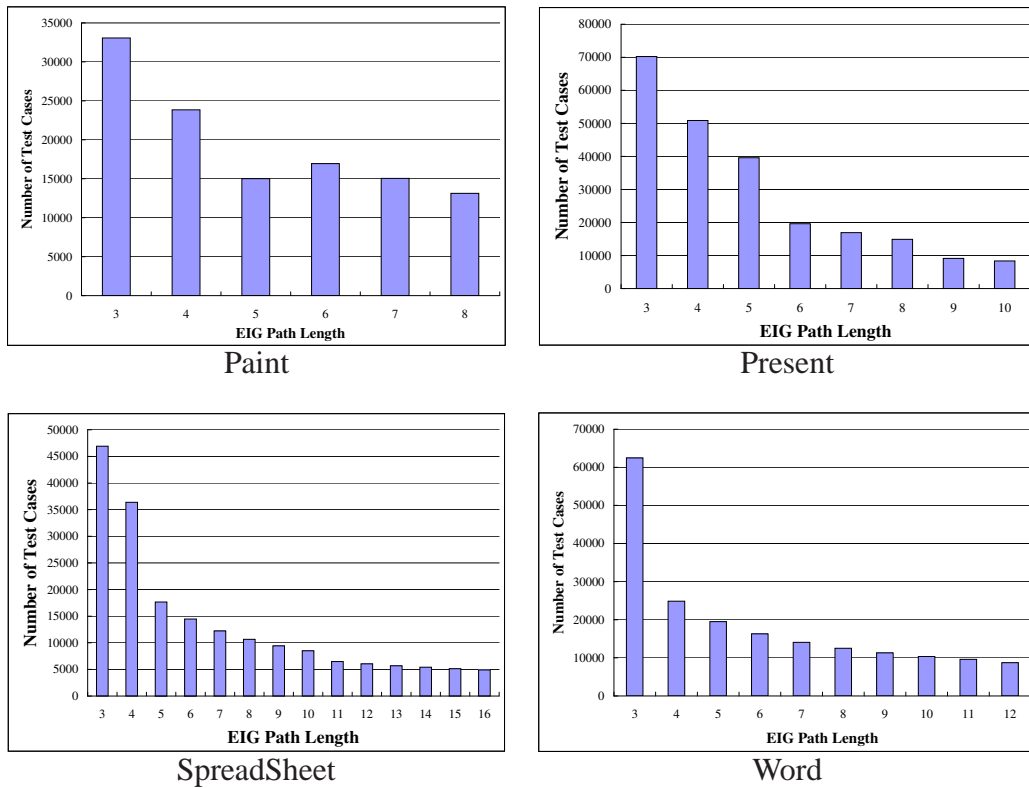


Figure 14: Histograms of Test Case Lengths in Pool

1 the entire process took over four months.

2 The total number of test cases per application is shown in Column 2 of Table 4. The length  
 3 distribution of the test cases is shown as a histogram in Figure 14. As expected, longer tests were  
 4 able to cover more EIG edges than the short ones; hence fewer long test cases were needed to  
 5 satisfy our coverage requirements.

6 After all the runs had completed, we had several matrices per application: (1) the fault matrix,  
 7 which summarized the faults detected per test case and (2) for each coverage criterion (event, edge,  
 8 statement), a coverage matrix, which summarized the coverage elements covered per test case.

9 This test pool was then used to obtain coverage-adequate suites. For example, event-adequate  
 10 suites were obtained by maintaining sets of test cases that covered each ESIG event. Test cases  
 11 were selected randomly without replacement from each set and duplicates eliminated, ensuring that  
 12 each event was covered by the resulting suite. A similar process was used for edge and statement



	Test Pool	$T_E$ Event	$T_I$ Edge	$T_S$ Stmt.	$T_R, T_{R+E}$ $T_{R+I}, T_{R+S}$
Paint	119,583	103	190	123.64	2531
Present	231,680	50	264	18.24	2080
SpreadSheet	191,966	45	173	14.08	2069
Word	192,042	84	248	30.35	2345

Table 4: Test Pool and Average-Suite Sizes

1 coverage. The process was repeated 100 times to yield 100 suites. The average size of the suites  
2 is shown in Columns 3–5 of Table 4.

3 Finally,  $T_R$  was constructed using random selection without replacement ensuring that the final  
4 size of  $T_R$  was the same as that of the ESIG suite. A total of 100 such suites per application were  
5 obtained. Similarly, each of the suites  $T_E$ ,  $T_I$ ,  $T_S$  were augmented with additional test cases,  
6 selected without replacement at random from the pool to yield  $T_{R+E}$ ,  $T_{R+I}$ ,  $T_{R+S}$ , respectively.  
7 The sizes of all these suites was equal to the size of the ESIG suite. Finally, 100 more suites that  
8 shared *all* the characteristics of interest in this study (*i.e.*, event, edge, statement, and size) with the  
9 ESIG suite were constructed; the symbol  $T_{R+E+I+S}$  will be used for these suites.

10 Note that the fault-detection effectiveness of each test suite can be obtained directly from the  
11 fault matrix of the test pool without rerunning the test cases. The results are shown in Figure 15  
12 as distributions. The box-plots provide a concise display of each distribution, each consisting  
13 of 100 data points. The line inside each box marks the median value. The edges of the box  
14 mark the first and third quartiles. The whiskers extend from the quartiles and cover 90% of the  
15 distribution; outliers are shown as points beyond the whiskers. Visual inspection of the plots  
16 shows that the fault-detection effectiveness of the ESIG-generated test suite (shown as an asterisk)  
17 is better than that of most individual similar-coverage and similar-sized suites. Some suites that  
18 lie in the whiskers and outliers do detect more faults than the ESIG suite. However, we remind  
19 the reader that unlike the ESIG suite, there is no systematic and automatic way to generate these  
20 suites.

21 As demonstrated above, box-plots are useful to get an overview of data distributions. However,  
22 valuable information is lost in creating the abstraction. For example, it is not clear *how many* test

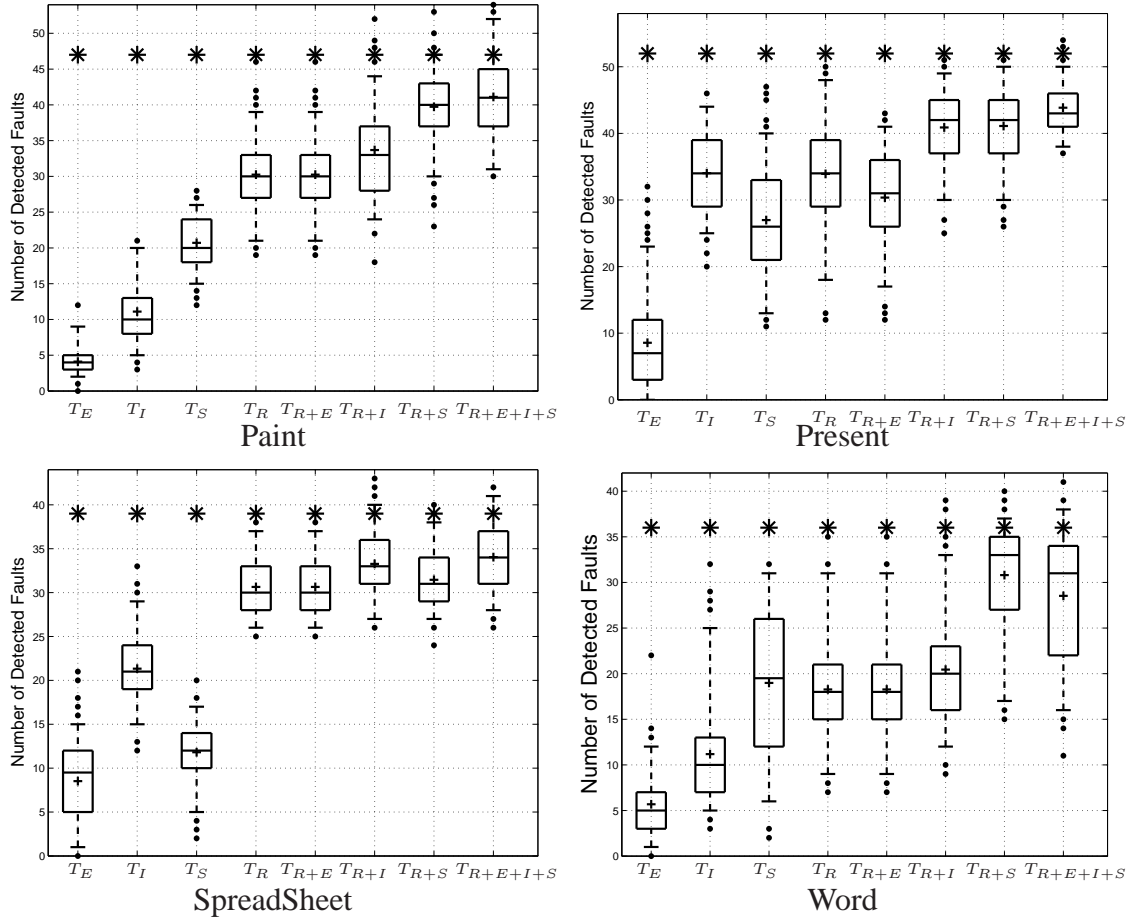


Figure 15: Fault Detection Distribution

1 suites detect specific numbers of faults; we are especially interested in the number of suites that do  
 2 better than the ESIG suites. This is important to partially understand our EIG and ESIG suites. We  
 3 now show the number of test suites that detected specific numbers of faults. Figure 17 shows eight  
 4 histograms for TerpPresent, one for each box in the box-plot. The x-axis represents the number of  
 5 faults; the y-axis shows the number of test suites that detected the particular number of faults. To  
 6 allow easy visual comparison, we have used the same x-axis and y-axis scales for all eight plots.  
 7 The vertical dotted line represents the number of faults detected by the ESIG suites. Similar plots  
 8 are shown in Figures 17–19 for Present, SpreadSheet, and Word, respectively.

9 For all applications,  $T_E$ ,  $T_I$ ,  $T_S$ ,  $T_R$ ,  $T_{R+E}$  consistently did worse than the ESIG suites. For  
 10 a small percentage of test suites,  $T_{R+I}$  and  $T_{R+S}$  did better than the ESIG suites. In particular,

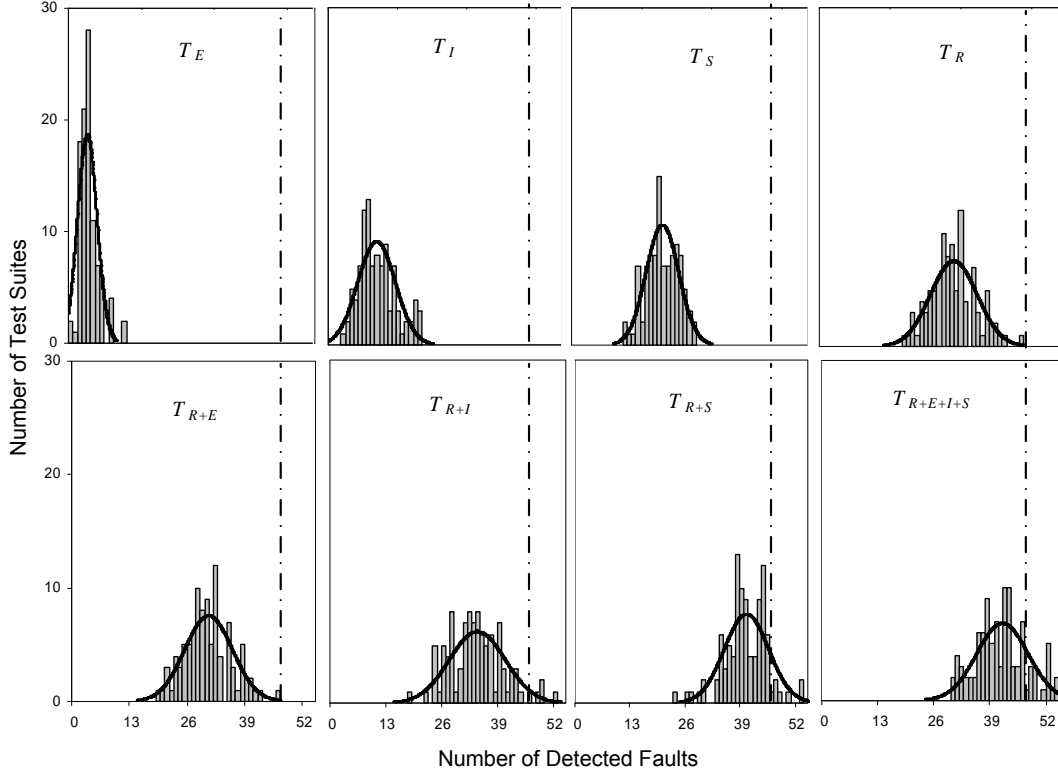


Figure 16: Histogram for Paint

1 4%, 0%, 6%, and 2% of the test suites for  $T_{R+I}$  did better than the ESIG suites, for Paint, Present,  
 2 SpreadSheet, and Word, respectively; 4%, 0%, 2%, and 11% of the test suites for  $T_{R+S}$  did better  
 3 than the ESIG suites, for Paint, Present, SpreadSheet, and Word, respectively.  $T_{R+E+I+S}$  outper-  
 4 formed all the other suites in most cases; Specifically, 15%, 2%, 8%, and 9% of the test suites for  
 5  $T_{R+E+I+S}$  did better than the ESIG suites, for Paint, Present, SpreadSheet, and Word, respectively.  
 6 It should be noted that the ESIG approach does better than most test suites considered in this study.  
 7 And that the ESIG-based approach is the only fully automatic approach to generate the GUI test  
 8 suites; all other suites were obtained from the test pool *after* they had been executed and statement  
 9 coverage obtained. Moreover, the performance of the  $T_R$  test suites indicates that the size of a suite  
 10 plays a very important role in fault-detection. We study this issue in the next section.

11 The results discussed thus far have been based on visual examination of the data. We now want  
 12 to determine whether the differences in the number of faults across coverage criteria are statistically  
 13 significant. In particular, we want to study the differences between the ESIG suite (a single value

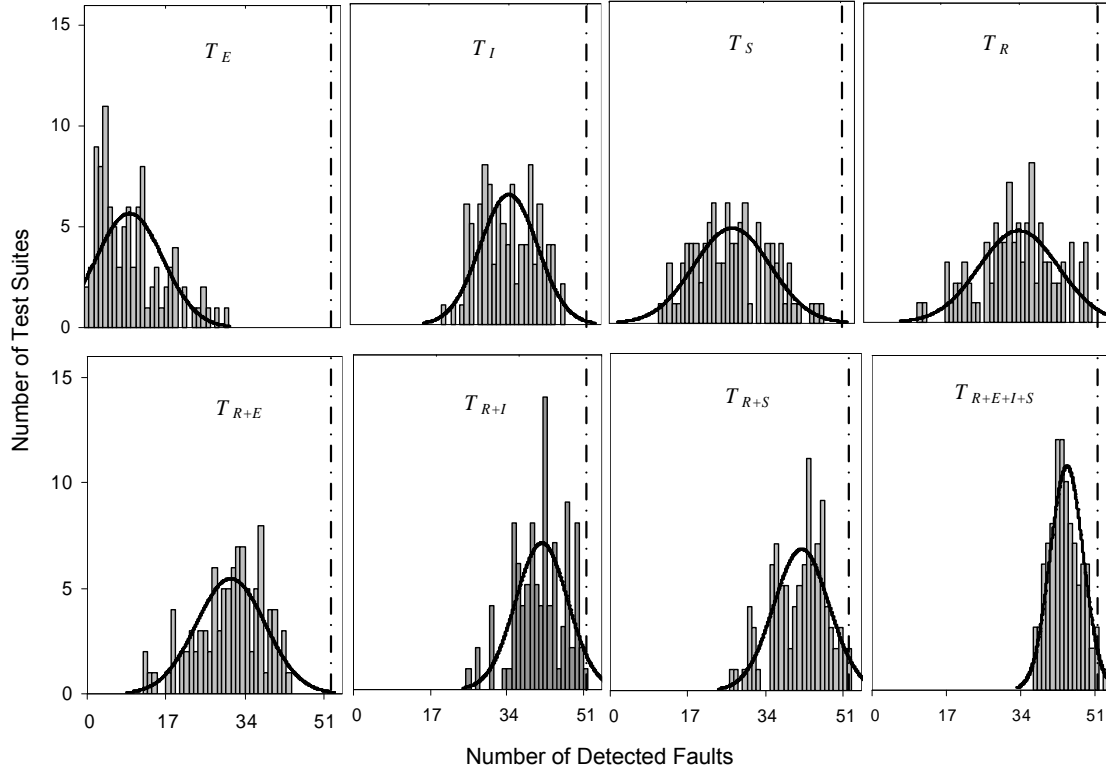


Figure 17: Histogram for Present

1 “number of faults detected”) and all the other “similar” suites (100 values per distribution). We  
 2 may use a *single sample t-test* to test for the significance of the difference between the observed  
 3 fault-detection of the ESIG suite and the mean of each distribution if certain conditions are met.

4 Many statistical tests are based upon the assumption that the data is normally distributed. These  
 5 tests are called *parametric* tests; a common example includes the t-test [55]. Since all the com-  
 6 putations in parametric tests are based on the data normality assumption, the significance level  
 7 (p-value) cannot be considered reliable when the data distribution is not normal. Tests that do not  
 8 make assumptions about distribution are referred to as *non-parametric* tests; these tests rank the  
 9 dependent variable from low to high and then analyze the ranks [55].

10 Another factor that plays a role in choice of statistical tests is sample size. If the sample  
 11 size is small, choosing a parametric test with non-normally distributed data, or choosing a non-  
 12 parametric test with normally distributed data, will yield an inaccurate p-value. On the other hand,  
 13 if the sample size is large, the distribution can be considered normal based upon the central limit

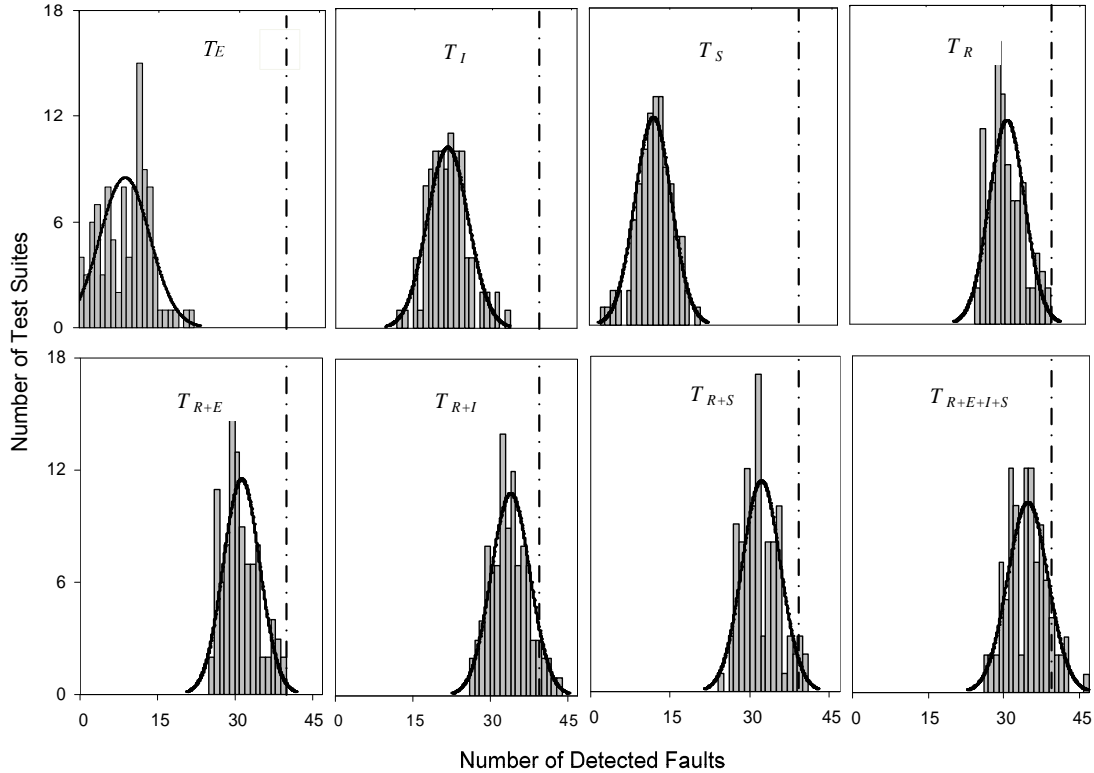


Figure 18: Histogram for SpreadSheet

1 theorem [55]. Consequently, in practice, with large sample size, the choice of parametric vs.  
 2 non-parametric does not matter, although the p-value tends to be large for non-parametric tests;  
 3 however, the discrepancy is small.

4 For illustration, the solid line superimposed on the histograms (Figures 17-19) shows the nor-  
 5 mal distribution approximation; informally, the data seems to follow the normal distribution quite  
 6 well. There are several ways to formally determine normality of data. A popular way is to use  
 7 a quantile-quantile (QQ) plot. A quantile is the fraction (or percentage) of points below a given  
 8 value. For example, 70% of the data fall below the “0.7 quantile” point. A QQ-plot shows the  
 9 quantiles of a data set (in our case the number of faults using different selection criteria) against  
 10 the quantiles of a second data set (normal distribution). A mark  $(x,y)$  on the plot means that  
 11  $f_2(x)$  and  $f_1(y)$  share the same quantile, where  $f_1$  and  $f_2$  are the functions for the first and second  
 12 data set. If the two sets come from a population with the same distribution, the plot will be along a  
 13  $(x = y)$  straight line. The greater the departure from this reference line, the greater the evidence for

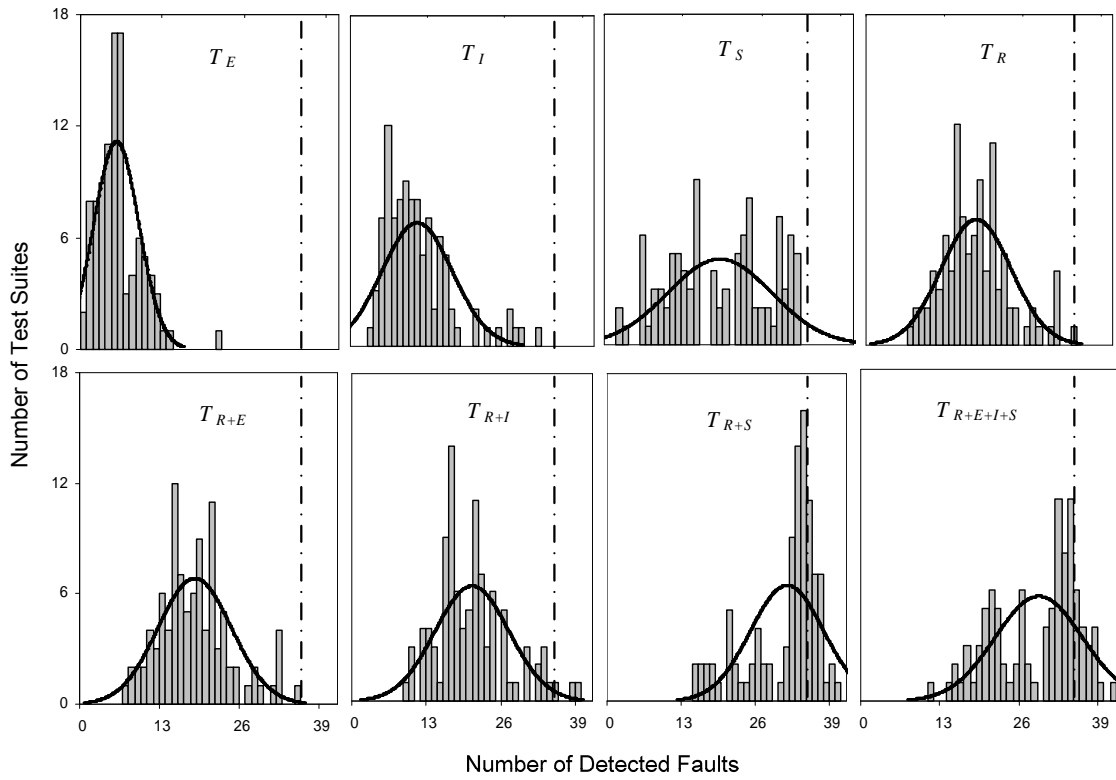


Figure 19: Histogram for Word

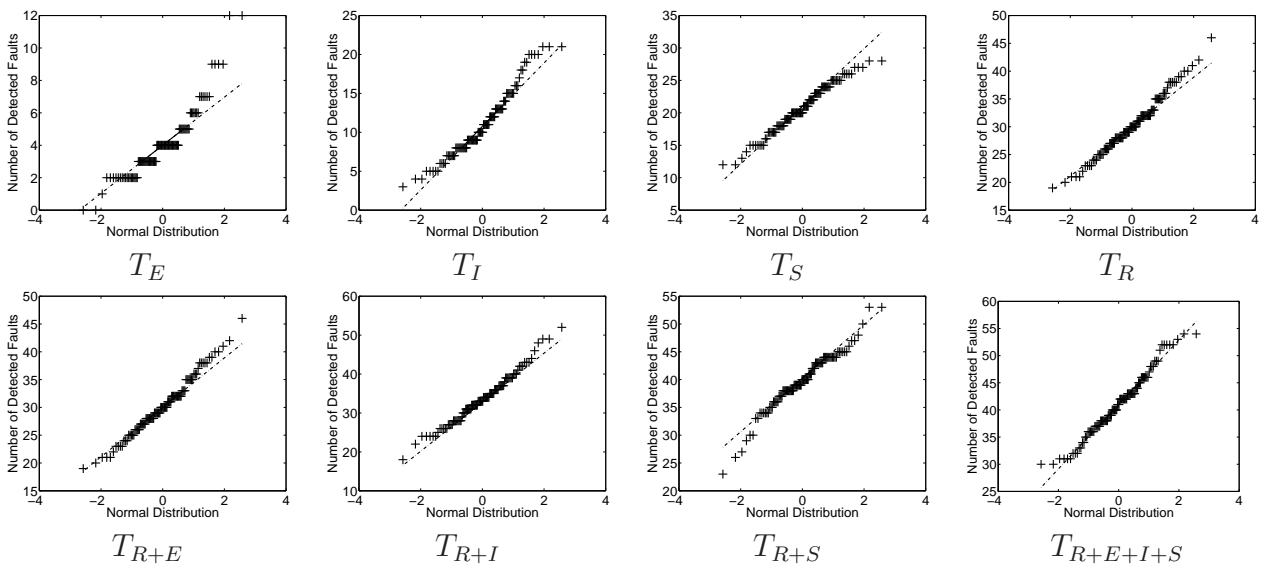


Figure 20: Quantile-Quantile plot for Paint

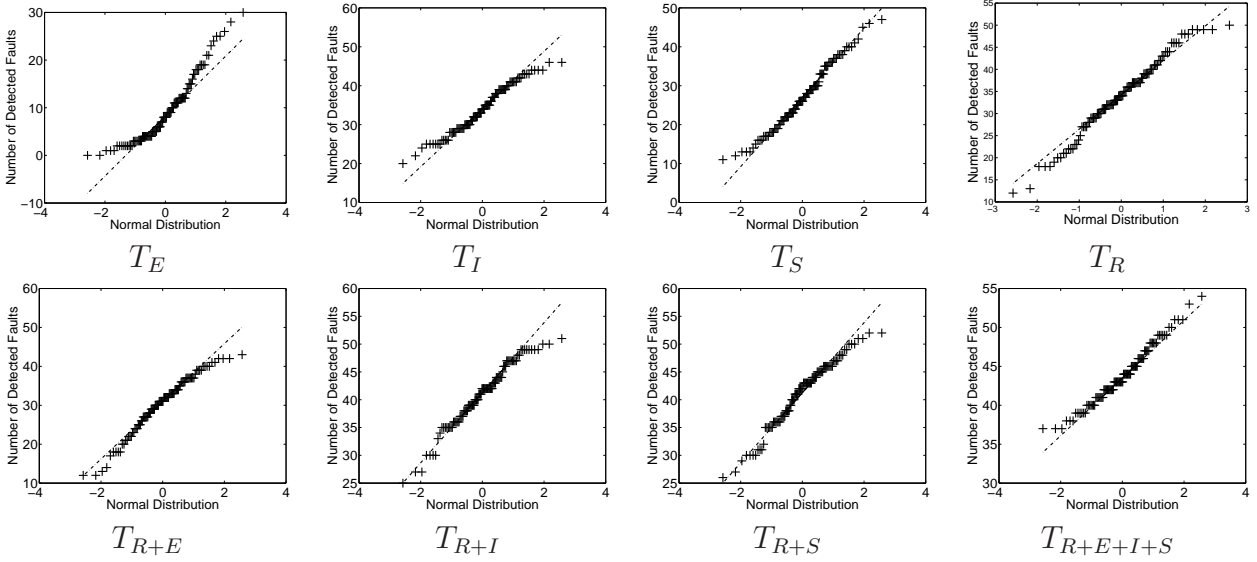


Figure 21: Quantile-Quantile plot for Present

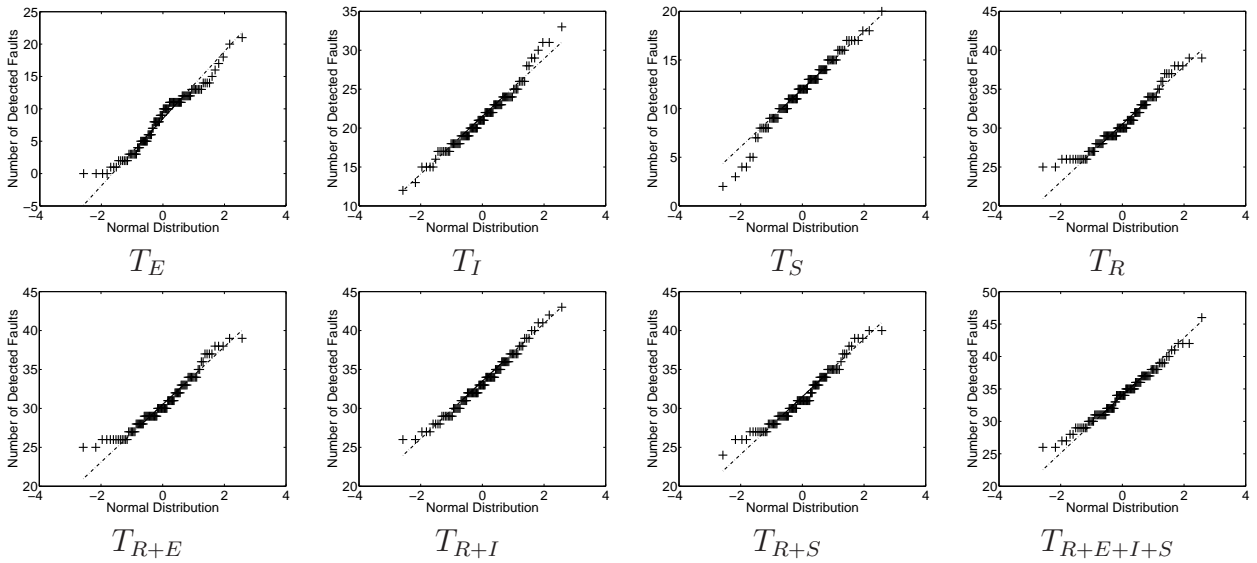


Figure 22: Quantile-Quantile plot for SpreadSheet

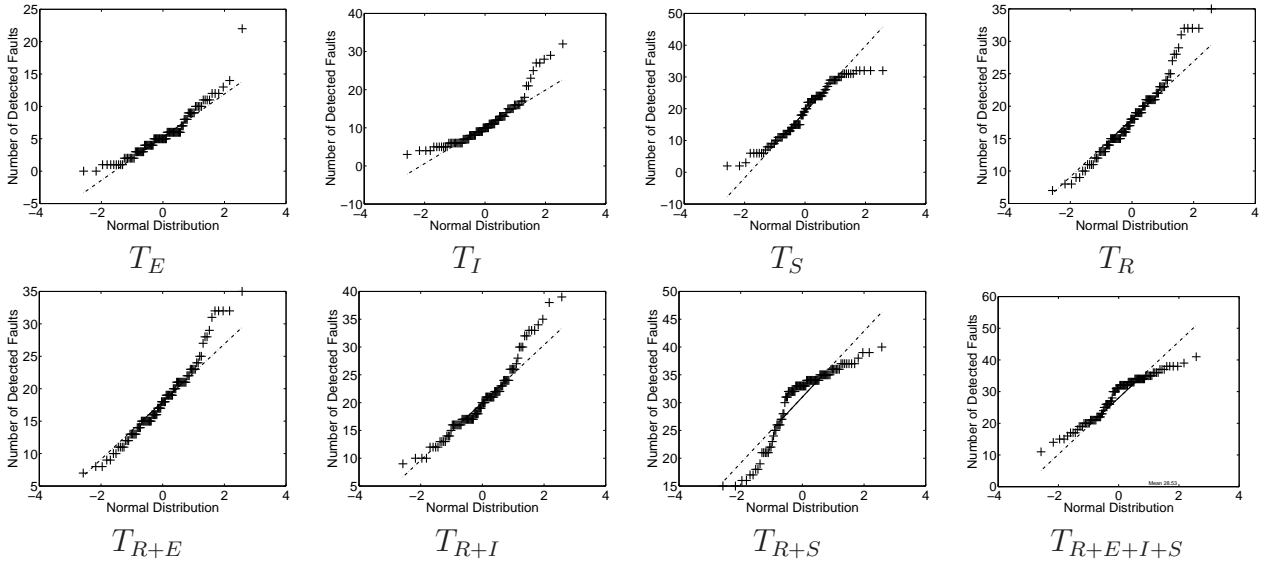


Figure 23: Quantile-Quantile plot for Word

1 the conclusion that the two data sets have come from populations with different distributions [60].  
 2 Figure 21 through Figure 22 show the quantile-quantile plots for each application. The x-axis is  
 3 the normal distribution; the y-axis is the “number of faults” distribution. The QQ-plots more or  
 4 less show a straight line, implying normality.

5 Based on normality, we use a *single sample t-test* to test for the significance of the difference  
 6 between the observed fault-detection of the ESIG suite and the mean of each distribution. The  
 7 null hypothesis is that the two values are equal; the alternate hypothesis is that they are unequal.  
 8 Note that a separate test is needed per pair (mean of each distribution, fault-detection of the ESIG  
 9 suite value). The resulting p-values were all more than 0.99. Hence we reject the null hypothesis  
 10 and accept the alternate hypothesis; there is a significant difference between fault-detection of the  
 11 ESIG suite and the mean fault-detection of each of the “similar” suites. The ESIG suites are better  
 12 at detecting faults compared to same-sized test suites that cover essentially the same events, edges,  
 13 and statements. This result helps to answer the primary question raised in this study.



## 6.4 Discussion

We now present details of why the ESIG-based test suites were able to detect more faults than other test suites. We specifically looked at three related issues: reachability, manifestation, and number of test cases. We note that the first two issues are related to the RELAY model [61, 70] of how a fault causes a failure on the execution of some test. We define *reachability* as the coverage of the statement in which a fault was seeded, and *manifestation* as the fault leading to a failure so that our GUI-based test oracle could detect it. As observed in [61], both are necessary conditions for fault detection. The data in Figure 15 indicates that the ESIG-based test suites were able to outperform their coverage-adequate equivalent counterparts. Hence, they must have been more successful, than their counterparts, at the combination of reachability and manifestation of several faults. We feel that this behavior is due to the nature of the ESI relationship, which is based on observed GUI state, and hence the software’s output. Executing tests that focus only on ESI events increases the likelihood that a fault will be revealed on the GUI, and hence detected by our GUI-based test oracle. Although we will not attempt to analyze this behavior in great detail here (indeed this is a direction for future research), we will provide some quantitative data to show evidence of this phenomenon by studying the test cases in the pool.

Figure 24 shows data of the number of test cases that cover a fault-seeded statement, and the number of test cases in the pool that detect it. The figure is divided into four parts, one for each subject application. Each part is a 3-column table (long tables are wrapped); Column 1 is the seeded fault number, corresponding to a statement in which the fault was seeded; Column 2 is the number of test cases that detected the fault; Column 3 is the number of test cases that did not detect the fault even though they executed the fault-seeded statement; the fault number entry is shaded if at least one ESI test case detected it. For example, the statement that contained Fault 21 of Paint was executed by a total of 845 test cases, of which 11 detected it; the fault was detected by at least one ESI test case. On the other hand, Fault 127 of Paint was not detected by any ESI test case; it was however detected by 42 of the total 42+267 test cases. Hence, a statement-coverage adequate test suite would have a probability of  $\frac{42}{42+267}$  of detecting this fault (assuming that faults



1 are independent). The data in Figure 24 is in fact sorted by this probability, giving us a sense  
2 of the “hardness” of the fault for statement-coverage adequate test suites. This data helps us to  
3 better interpret the results of Figure 15. First, the ESIG test suites did detect many of the seeded  
4 faults. Second, they did better than  $T_S$  because they detected many of the “hard” faults (this was  
5 most apparent in Paint and SpreadSheet). Third, some faults were detected by many of the test  
6 cases that executed the statement. For example, Fault 136 in Paint was detected by 747 of the total  
7 747+446 test cases that executed the statement in which it was seeded. The fault was seeded in the  
8 handler of a termination event that closes the `Attributes` window and applies the attributes (if  
9 any have changed) to the current image on the main canvas. The seeded fault caused the image size  
10 to be computed incorrectly, resulting in an incorrectly sized image whenever at least one attribute  
11 in the `Attributes` window is modified by the user. Because there are many permutations of  
12 modifying the attributes, a large number of test cases are able to detect this fault (12 in the ESIG  
13 test suite and 735 in the rest of the test pool). In general, statement coverage adequate test suites  
14 do really well for these types of faults that can be triggered in many different ways.

15 Forth, several faults were detected by very few test cases. Fault 34 of Paint is an example. This  
16 fault flips the conditional statement in an event handler of a type of curve tool. The condition is  
17 to check whether the curve tool is currently selected; if yes, then the curve stroke is set according  
18 to the selected line type for the curve tool. Due to the fault, the curve stroke is incorrectly set,  
19 resulting in an incorrect image to be drawn on the main canvas only when the event sequence: `<`  
20 `SelectCurveTool; SelectLineType; DrawOnCanvas >` is executed. If the first two events are  
21 not executed, then `DrawOnCanvas` does not trigger the fault even if the statement containing the  
22 seeded fault is executed. Hence, although there are many test cases that cover the faulty statement  
23 (850), only 7 test cases in the test pool detected the fault. One of them is the ESI test case; ESI-  
24 relationships were found between the three events. In general, statement coverage adequate test  
25 suites do not do well for faults that are executed by many event sequences but manifest as failures  
26 in very few cases.

27 The size of a suite seems to play a very important role in fault-detection. Indeed, the  $T_R$  test



1 suites, which were the same size as the ESIG-based suites, did better (in most cases) than their  
2 coverage-adequate counterparts. We feel that this behavior is an artifact of the density of our fault  
3 matrices. A large number of test cases are successful at detecting many “easy” faults. Even if test  
4 cases are selected at random, given adequate numbers, they will be able to detect a large number  
5 of these easy faults. For example, given 192,042 test cases in the pool for Word and the size of  $T_R$   
6 being 2345, the probability that Fault 24, which is detected by more than 84 test cases in the pool,  
7 is detected by at least one test case in  $T_S$  is 0.49 – this is quite high. In Figure 25, we show the  
8 probability that a random suite of its corresponding ESI-suite size would detect a particular fault.  
9 This data shows that many of these difficult faults are detected by at least one ESIG-based test  
10 case, improving their fault-detection effectiveness. Moreover, 16 faults in Word have a detection  
11 probability of more than 0.25. This number is much larger for the other three applications, helping  
12 to understand why  $T_R$  and the other suites that included randomly selected test cases did so well.

13 Finally, we wanted to examine why some of the faults were not detected. We manually ex-  
14 amined each fault and tried to manually devise ways of manifesting the fault as a failure. We  
15 determined that:

- 16 1. several of the faults were in fact manifested as failures on the GUI but our test oracle was  
17 not capable of examining these parts of the GUI,
- 18 2. very few faults caused failures in non-GUI output,
- 19 3. several of the undetected faults require even longer sequences,
- 20 4. the effect of several faults was masked by the event handler code even before our test oracle  
21 could detect it,
- 22 5. some faults crashed their corresponding fault-seeded version.

23 We show the numbers of these faults in Table 5. The large number of “Ignored Widget Proper-  
24 ties” has prompted us to improve our test oracles for future work.

25 This controlled study showed that the automatically identified ESI relationships between events  
26 generate test suites that detect more faults than their code-, event-, and event-interaction-coverage  
27 equivalent counterparts. Moreover, we saw that several of our missed faults remained undetected

	Ignored Widget Properties	Non-GUI Failure	Longer Sequence	Masked Error	Crash	Total
Paint	0	0	1	6	0	7
Present	13	0	4	0	0	17
SpreadSheet	9	2	8	5	3	27
Word	5	0	4	11	0	20

Table 5: Undetected Faults Classification

1 because of limitations with our automated GUI-based test oracle, and others required even longer  
2 sequences.

## 3 7 Conclusions and Future Work

4 This paper presented a new fully automatic technique to test multi-way interactions among GUI  
5 events. The technique is based on analysis of feedback obtained from the run-time state of GUI  
6 widgets. A seed test suite is used for feedback collection. The technique was demonstrated via two  
7 independent studies on eight software applications. The results of the first study showed that the  
8 test cases generated using the feedback were useful at detecting serious and relevant faults in the  
9 applications. The second study compared the ESIG-based test suite to similar EIG-based suites. It  
10 showed that the added effectiveness is due to targeted testing of the identified ESI relationships, not  
11 an incidental side-effect of the size of the suite, nor the additional events and code that it covers.

12 As is the case with all research involving empirical studies, these studies are subject to threats  
13 to validity. These threats need to be considered in order to assess their impact on the results.  
14 First is the selection of subject applications and their characteristics. The results may vary for  
15 applications that have a complex back-end, are not developed using the object-oriented paradigm,  
16 or have non-deterministic behavior. Second, in study 2, the test pool approach was used due  
17 to practical limitations. It is expected that the repetition of the same test case across multiple test  
18 suites will have an impact on some of the results. The algorithm used to create the test pool ensures  
19 that each event (the first event in the test case) is executed in a known initial state; the choice of  
20 this state may have an effect on the results. Third, the Java API allow the extraction of only 12

1 properties of each widget; only these properties were used for obtaining the ESI relationship via  
2 GUI state; moreover, faults are reported for mismatches between these 12 properties. Fourth,  
3 we used one technique to generate test cases – using event-interaction graphs. Other techniques,  
4 *e.g.*, using capture/replay tools and programming the test cases manually may produce different  
5 types of test cases, which may show different execution behavior. Fifth, in study 2, a threat is  
6 related to our measurement of fault detection effectiveness; each fault was seeded and activated  
7 individually. Note that multiple faults present simultaneously can lead to more complex scenarios  
8 that include fault masking. Finally, several threats are related to fault seeding in study 2. Threats  
9 from issues such as human decision-making are minimized by using an objective technique for  
10 uniformly distributing faults based on functional units.

11 This research has presented several exciting opportunities for future work. In the immediate  
12 future, the three contexts for the cases will be simplified and, if possible, combined. The current  
13 special treatment of termination events, which led to an additional two contexts, will be revised.  
14 One possibility is the revision of the EIG model; the elimination of all termination events from  
15 this model will be explored. This revision will also lead to the definition of new, fundamentally  
16 different cases for the ESI relationship.

17 The results showed that certain events in the GUI dominate the ESI relationship. These events  
18 will be studied and classified. In the future, additional GUI applications and software problems  
19 will be studied. The run-time state information was collected using the Java Swing API for stan-  
20 dard Swing widgets. Future work involves incorporating customized API for application-specific  
21 widgets into feedback collection and analysis.

22 The current test-case generation algorithms output a set of all possible event sequences bounded  
23 by a pre-determined length. As the goal of this work is to generate multi-way interactions among  
24 GUI events, other techniques (such as *covering arrays* [85]) designed to minimize the number of  
25 test cases while retaining high interaction coverage will be explored.

26 The analysis summarized in Section 5 led to a deeper understanding of the relationship be-  
27 tween real GUI events and the underlying code in fielded GUI applications. This may lead to

1 new techniques that combine dynamic analysis of the GUI and static analysis of the event handler  
2 code. For example, the code for related events may be given to a static-analysis engine that could  
3 examine the code for possible interactions that are only apparent at the code level, *e.g.*, data-flow  
4 relationships.

5 The feedback currently obtained at run time is in the form of GUI widgets. Mechanisms, such  
6 as reflection, in modern programming languages may be used to obtain additional feedback from  
7 non-GUI objects. The definition of state, in terms of a set of objects with properties and values, is  
8 general; it may be applied to any executing object. Some of the six cases may be adapted for non-  
9 GUI objects. Another straightforward way to enhance the feedback is to instrument the software  
10 for code coverage and run-time invariant collection. This feedback may be used to generate new  
11 types of test cases. Another logical extension of this work is to examine the redundancy in our  
12 ESIG suites via existing test minimization techniques developed for user interfaces [9].

13 Some of the challenges of GUI testing are also relevant to testing of event-driven software,  
14 *e.g.*, web applications and object-oriented software. One way to test these classes of software is  
15 to generate test cases that are sequences of events (*e.g.*, web user actions or method calls). Some  
16 of the techniques developed in this research have already been used by other researchers to prune  
17 the space of all possible event interactions to be tested for web applications [2]; similar extensions  
18 will be explored for object-oriented software.

## 19 **Acknowledgments**

20 We thank the anonymous reviewers whose comments and suggestions helped to extend the second  
21 empirical study, reshape its results, and improve the flow of the text. This work was partially  
22 supported by the US National Science Foundation under NSF grant CCF-0447864 and the Office  
23 of Naval Research grant N00014-05-1-0421.



## 1 **References**

- 2 [1] AHO, A. V., DAHBURA, A. T., LEE, D., AND UYAR, M. U. An optimization technique  
3 for protocol conformance test generation based on uio sequences and rural chinese postman  
4 tours. *Conformance testing methodologies and architectures for OSI protocols* (1995), 427–  
5 438.
- 6 [2] ALESSANDRO MARCHETTO, P. T., AND RICCA, F. State-based testing of Ajax web appli-  
7 cations. In *Proceedings of the 1st International Conference on Software Testing, Verification,*  
8 *and Valication* (April 9–11, 2008), pp. 121–130.
- 9 [3] ALLEN, F. E. Control flow analysis. In *Proceedings of a symposium on Compiler optimiza-*  
10 *tion* (1970), pp. 1–19.
- 11 [4] APFELBAUM, L. Automated functional test generation. In *Autotestcon '95 Conference*  
12 (1995), IEEE.
- 13 [5] AUGUSTON, M., MICHAEL, J. B., AND SHING, M.-T. Environment behavior models for  
14 scenario generation and testing automation. In *A-MOST '05: Proceedings of the 1st inter-*  
15 *national workshop on Advances in model-based testing* (New York, NY, USA, 2005), ACM  
16 Press, pp. 1–6.
- 17 [6] BARNETT, M., GRIESKAMP, W., NACHMANSON, L., SCHULTE, W., TILLMANN, N., AND  
18 VEANES, M. Towards a tool environment for model-based testing with AsmL. In Petrenko  
19 and Ulrich [59], pp. 252–266.
- 20 [7] BELLI, F. Finite-state testing and analysis of graphical user interfaces. In *ISSRE* (2001),  
21 IEEE Computer Society, pp. 34–43.
- 22 [8] BELLI, F. Finite-state testing and analysis of GUIs. In *Proceedings of the 12th International*  
23 *Symposium on Software Reliability Engineering* (2001), pp. 34–43.

- 1 [9] BELLI, F., AND BUDNIK, C. J. Test minimization for human-computer interaction. *Applied*  
2 *Intelligence* 26, 2 (2007), 161–174.
- 3 [10] BELLI, F., BUDNIK, C. J., AND WHITE, L. Event-based modelling, analysis and testing of  
4 user interactions: approach and case study: Research articles. *Softw. Test. Verif. Reliab.* 16, 1  
5 (2006), 3–32.
- 6 [11] BERNHARD, P. J. A reduced test suite for protocol conformance testing. *ACM Transactions*  
7 *on Software Engineering and Methodology* 3, 3 (July 1994), 201–220.
- 8 [12] BLACKBURN, M. R., BUSSEY, R., NAUMAN, A., AND CHANDRAMOULI, R. Model-  
9 based approach to security test automation. In *13th International Symposium on Software*  
10 *Reliability Engineering (ISSRE)* (Annapolis, MD, 2002).
- 11 [13] BLACKBURN, M. R., BUSSEY, R., NAUMAN, A., KNICKERBOCKER, R., AND KASUDA,  
12 R. Mars polar lander fault identification using model-based testing. In *ICECCS (2002)*, IEEE  
13 Computer Society, pp. 163–.
- 14 [14] BOYAPATI, C., KHURSHID, S., AND MARINOV, D. Korat: automated testing based on java  
15 predicates. In *ISSTA '02: Proceedings of the 2002 ACM SIGSOFT international symposium*  
16 *on Software testing and analysis* (2002), pp. 123–133.
- 17 [15] BRIAND, L. C., LABICHE, Y., AND WANG, Y. Using simulation to empirically investigate  
18 test coverage criteria based on statechart. In *ICSE '04: Proceedings of the 26th International*  
19 *Conference on Software Engineering* (2004), IEEE Computer Society, pp. 86–95.
- 20 [16] CAMPBELL, C., GRIESKAMP, W., NACHMANSON, L., SCHULTE, W., TILLMANN, N.,  
21 AND VEANES, M. Model-based testing of object-oriented reactive systems with spec ex-  
22 plorer., May 2005.

- 1 [17] CHEN, W.-H., LU, C.-S., BROZOVSKY, E. R., AND WANG, J.-T. An optimization tech-  
2 nique for protocol conformance testing using multiple uio sequences. *Inf. Process. Lett.* 36,  
3 1 (1990), 7–11.
- 4 [18] CHOW, T. S. Testing software design modeled by finite-state machines. *IEEE trans. on*  
5 *Software Engineering SE-4*, 3 (1978), 178–187.
- 6 [19] CHOW, T. S. Testing software design modeled by finite-state machines. *IEEE Trans. Softw.*  
7 *Eng.* 4, 3 (1978), 178–187.
- 8 [20] CLARKE, J. M. Automated test generation from a behavioral model. In *Proceedings of*  
9 *Pacific Northwest Software Quality Conference* (Portland, OR, May 1998), Pnsqc/Pacific  
10 Agenda.
- 11 [21] D’AMORIM, M., PACHECO, C., XIE, T., MARINOV, D., AND ERNST, M. D. An em-  
12 pirical comparison of automated generation and classification techniques for object-oriented  
13 unit testing. In *Proceedings of the 21st IEEE/ACM International Conference on Automated*  
14 *Software Engineering* (2006).
- 15 [22] DUSTIN, E., RASHKA, J., AND PAUL, J. *Automated software testing: introduction, man-*  
16 *agement, and performance.* 1999.
- 17 [23] ESMELIOGLU, S., AND APFELBAUM, L. Automated test generation, execution, and report-  
18 ing. In *Proceedings of Pacific Northwest Software Quality Conference* (Portland, OR, Oct  
19 1997), Pnsqc/Pacific Agenda, pp. 127–142.
- 20 [24] FARCHI, E., HARTMAN, A., AND PINTER, S. S. Using a model-based test generator to test  
21 for standard conformance. *IBM Systems Journal* 41, 1 (2002), 89–110.
- 22 [25] FERGUSON, R., AND KOREL, B. The chaining approach for software test data generation.  
23 *ACM Trans. Softw. Eng. Methodol.* 5, 1 (1996), 63–86.

- 1 [26] GALLAGHER, M. J., AND NARASIMHAN, V. L. Adtest: A test data generation suite for  
2 Ada software systems. *IEEE Trans. Software Eng.* 23, 8 (1997), 473–484.
- 3 [27] GUPTA, N., MATHUR, A. P., AND SOFFA, M. L. Automated test data generation using an  
4 iterative relaxation method. In *SIGSOFT FSE* (1998), pp. 231–244.
- 5 [28] HONG, H. S., KWON, Y. R., AND CHA, S. D. Testing of object-oriented programs based  
6 on finite state machines. In *APSEC* (1995), IEEE Computer Society, pp. 234–.
- 7 [29] HOVEMEYER, D., AND PUGH, W. Finding bugs is easy. *SIGPLAN Not.* 39, 12 (2004),  
8 92–106.
- 9 [30] HOWE, A., VON MAYRHAUSER, A., AND MRAZ, R. T. Test case generation as an AI  
10 planning problem. *Automated Software Engineering* 4 (1997), 77–106.
- 11 [31] IMANIAN, J. A. *Automatic Test Case Generation for Reactive Software Systems Based on*  
12 *Environment Models*. PhD thesis, Naval Postgraduate School, Monterey, CA, June 2005.
- 13 [32] IPATE, F., AND HOLCOMBE, M. Specification and testing using generalized machines: a  
14 presentation and a case study. *Software Testing, Verification and Reliability* (1998), 61–81.
- 15 [33] IPATE, F., AND HOLCOMBE, M. Complete testing from a stream x-machine specification.  
16 *Fundam. Inf.* 64, 1-4 (2004), 205–216.
- 17 [34] JORGENSEN, A. A., AND WHITTAKER, J. A. An application program interface (API) test-  
18 ing method. In *STAREast Conference* (May 2000).
- 19 [35] KASIK, D. J., AND GEORGE, H. G. Toward automatic generation of novice user test scripts.  
20 In *Proceedings of the Conference on Human Factors in Computing Systems : Common*  
21 *Ground* (New York, 13–18 Apr. 1996), ACM Press, pp. 244–251.
- 22 [36] KOOPMAN, P. W. M., PLASMEIJER, R., AND ACHTEN, P. Model-based testing of thin-  
23 client web applications. In *FATES/RV* (2006), K. Havelund, M. Núñez, G. Rosu, and  
24 B. Wolff, Eds., vol. 4262 of *Lecture Notes in Computer Science*, Springer, pp. 115–132.

- 1 [37] KOREL, B. Automated software test data generation. *IEEE Trans. Software Eng.* 16, 8  
2 (1990), 870–879.
- 3 [38] LEE, N. H., AND CHA, S. D. Generating test sequences from a set of mscs. *Computer*  
4 *Networks* 42, 3 (2003), 405–417.
- 5 [39] LEOW, W. K., KHOO, S. C., AND SUN, Y. Automated generation of test programs from  
6 closed specifications of classes and test cases. In *ICSE '04: Proceedings of the 26th Interna-*  
7 *tional Conference on Software Engineering* (Washington, DC, USA, 2004), IEEE Computer  
8 Society, pp. 96–105.
- 9 [40] LUCIO, L., PEDRO, L., AND BUCHS, D. A methodology and a framework for model-based  
10 testing. In *RISE (2004)*, N. Guelfi, Ed., vol. 3475 of *Lecture Notes in Computer Science*,  
11 Springer, pp. 57–70.
- 12 [41] MATHUR, A. P. *Foundations of Software Testing: Fundamental Algorithms and Techniques*.  
13 Pearson Education., 2008.
- 14 [42] MAURER, P. M. Generating test data with enhanced context-free grammars. *IEEE Software*  
15 7, 4 (July 1990), 50–55.
- 16 [43] MCMMASTER, S., AND MEMON, A. M. Call-stack coverage for GUI test-suite reduction.  
17 *IEEE Trans. Softw. Eng.* (2008).
- 18 [44] MEMON, A. M. An event-flow model of gui-based applications for testing. *Software Testing,*  
19 *Verification and Reliability* 17, 3 (2007), 137–157.
- 20 [45] MEMON, A. M., NAGARAJAN, A., AND XIE, Q. Automating regression testing for evolving  
21 GUI software. *Journal of Software Maintenance and Evolution* 17, 1 (Jan. 2005), 27–64.
- 22 [46] MEMON, A. M., NAGARAJAN, A., AND XIE, Q. Automating regression testing for evolving  
23 gui software. *Journal of Software Maintenance* 17, 1 (2005), 27–64.

- 1 [47] MEMON, A. M., POLLACK, M. E., AND SOFFA, M. L. Hierarchical GUI test case gen-  
2 eration using automated planning. *IEEE Transactions on Software Engineering* 27, 2 (Feb.  
3 2001), 144–155.
- 4 [48] MEMON, A. M., POLLACK, M. E., AND SOFFA, M. L. Hierarchical GUI test case genera-  
5 tion using automated planning. *IEEE Trans. Softw. Eng.* 27, 2 (2001), 144–155.
- 6 [49] MEMON, A. M., AND XIE, Q. Studying the fault-detection effectiveness of GUI test cases  
7 for rapidly evolving software. *IEEE Trans. Softw. Eng.* 31, 10 (2005), 884–896.
- 8 [50] MEMON, A. M., AND XIE, Q. Studying the fault-detection effectiveness of GUI test cases  
9 for rapidly evolving software. *IEEE Transactions on Software Engineering* 31, 10 (Oct.  
10 2005), 884–896.
- 11 [51] MEMON, A. M., AND XIE, Q. Studying the fault-detection effectiveness of GUI test cases  
12 for rapidly evolving software. *IEEE Trans. Softw. Eng.* 31, 10 (2005), 884–896.
- 13 [52] MICHAEL, C. C., MCGRAW, G., AND SCHATZ, M. Generating software test data by evo-  
14 lution. *IEEE Trans. Software Eng.* 27, 12 (2001), 1085–1110.
- 15 [53] MICSKEI, Z., AND MAJZIK, I. Model-based automatic test generation for event-driven  
16 embedded systems using model checkers. In *DepCoS-RELCOMEX* (2006), IEEE Computer  
17 Society, pp. 191–198.
- 18 [54] MILLER, W., AND SPOONER, D. L. Automatic generation of floating-point test data. *IEEE*  
19 *Trans. Software Eng.* 2, 3 (1976), 223–226.
- 20 [55] MOTULSKY, H. *Intuitive Biostatistics*. Oxford University Press, 1995.
- 21 [56] MYERS, B. A., AND ROSSON, M. B. Survey on user interface programming. In *CHI* (1992),  
22 pp. 195–202.

- 1 [57] PACHECO, C., LAHIRI, S. K., ERNST, M. D., AND BALL, T. Feedback-directed random  
2 test generation. In *ICSE '07: Proceedings of the 29th International Conference on Software*  
3 *Engineering* (Minneapolis, MN, USA, May 23–25, 2007), IEEE Computer Society, pp. 396–  
4 405.
- 5 [58] PACHECO, C., LAHIRI, S. K., ERNST, M. D., AND BALL, T. Feedback-directed random  
6 test generation. In *ICSE'07, Proceedings of the 29th International Conference on Software*  
7 *Engineering* (Minneapolis, MN, USA, May 23–25, 2007).
- 8 [59] PETRENKO, A., AND ULRICH, A., Eds. *Formal Approaches to Software Testing, Third In-*  
9 *ternational Workshop on Formal Approaches to Testing of Software, FATES 2003, Montreal,*  
10 *Quebec, Canada, October 6th, 2003* (2004), vol. 2931 of *Lecture Notes in Computer Science*,  
11 Springer.
- 12 [60] Quantile-Quantile Plot. <http://www.itl.nist.gov/div898/handbook/eda/section3/qqplot.htm> .
- 13 [61] RICHARDSON, D. J., AND THOMPSON, M. C. An analysis of test data selection criteria  
14 using the relay model of fault detection. *IEEE Trans. Softw. Eng.* 19, 6 (1993), 533–553.
- 15 [62] ROSEN, B. K. Data flow analysis for procedural languages. *J. ACM* 26, 2 (1979), 322–344.
- 16 [63] ROTHERMEL, G., ELBAUM, S., MALISHEVSKY, A. G., KALLAKURI, P., AND QIU, X.  
17 On test suite composition and cost-effective regression testing. *ACM Trans. Softw. Eng.*  
18 *Methodol.* 13, 3 (2004), 277–331.
- 19 [64] ROUNTEV, A., KAGAN, S., AND GIBAS, M. Evaluating the imprecision of static analysis.  
20 In *Proceedings of the ACM-SIGPLAN-SIGSOFT workshop on Program analysis for software*  
21 *tools and engineering* (2004), pp. 14–16.
- 22 [65] SARIKAYA, B. Conformance testing: architectures and test sequences. *Comput. Netw. ISDN*  
23 *Syst.* 17, 2 (1989), 111–126.

- 1 [66] SCHEETZ, M., VON MAYRHAUSER, A., FRANCE, R., DAHLMAN, E., AND HOWE, A. E.  
2 Generating test cases from an OO model with an ai planning system. In *Proceedings of The*  
3 *10th International Symposium on Software Reliability Engineering* (Washington, DC, USA,  
4 1999), IEEE Computer Society, pp. 250–259.
- 5 [67] SHEHADY, R. K., AND SIEWIOREK, D. P. A method to automate user interface testing  
6 using variable finite state machines. In *FTCS '97: Proceedings of the 27th International*  
7 *Symposium on Fault -Tolerant Computing (FTCS '97)* (Washington, DC, USA, 1997), IEEE  
8 Computer Society, p. 80.
- 9 [68] SHEHADY, R. K., AND SIEWIOREK, D. P. A method to automate user interface testing using  
10 variable finite state machines. In *Proceedings of The Twenty-Seventh Annual International*  
11 *Symposium on Fault-Tolerant Computing (FTCS'97)* (Washington - Brussels - Tokyo, June  
12 1997), IEEE Press, pp. 80–88.
- 13 [69] SIRER, E. G., AND BERSHAD, B. N. Using production grammars in software testing. In  
14 *PLAN '99: Proceedings of the 2nd conference on Domain-specific languages* (New York,  
15 NY, USA, 1999), ACM Press, pp. 1–13.
- 16 [70] THOMPSON, M. C., RICHARDSON, D. J., AND CLARKE, L. A. An information flow model  
17 of fault detection. In *ISSTA '93: Proceedings of the 1993 ACM SIGSOFT international*  
18 *symposium on Software testing and analysis* (New York, NY, USA, 1993), ACM, pp. 182–  
19 192.
- 20 [71] VON MAYRHAUSER, A., AND CRAWFORD-HINES, S. Automated testing support for a robot  
21 tape library. In *Proceedings of The Fourth International Software Reliability Engineering*  
22 *Conference* (Washington, DC, USA, Nov. 1993), IEEE Computer Society, pp. 6–14.
- 23 [72] VON MAYRHAUSER, A., MRAZ, R. T., AND WALLS, J. Domain based regression testing.  
24 In *Proceedings of The International Conference on Software Maintenance* (Washington, DC,  
25 USA, 1994), IEEE Computer Society, pp. 26–35.



- 1 [73] WHITE, L., AND ALMEZEN, H. Generating test cases for GUI responsibilities using com-  
2 plete interaction sequences. In *ISSRE '00: Proceedings of the 11th International Symposium*  
3 *on Software Reliability Engineering* (Washington, DC, USA, 2000), IEEE Computer Soci-  
4 ety, p. 110.
- 5 [74] WHITE, L., AND ALMEZEN, H. Generating test cases for gui responsibilities using complete  
6 interaction sequences. In *ISSRE '00: Proceedings of the 11th International Symposium on*  
7 *Software Reliability Engineering (ISSRE'00)* (Washington, DC, USA, 2000), IEEE Computer  
8 Society, p. 110.
- 9 [75] WHITTAKER, J. A. *Markov chain techniques for software testing and reliability analysis*.  
10 PhD thesis, University of Tennessee, Knoxville, TN, USA, 1992.
- 11 [76] WHITTAKER, J. A. Stochastic software testing. *Ann. Software Eng.* 4 (1997), 115–131.
- 12 [77] WHITTAKER, J. A., AND THOMASON, M. G. A markov chain model for statistical software  
13 testing. *IEEE Trans. Softw. Eng.* 20, 10 (1994), 812–824.
- 14 [78] WOIT, D. Conditional-event usage testing. In *CASCON '98: Proceedings of the 1998 con-*  
15 *ference of the Centre for Advanced Studies on Collaborative research* (Indianapolis, Indiana,  
16 USA, 1998), IBM Press, p. 23.
- 17 [79] WOIT, D. M. Specifying operational profiles for modules. In *ISSTA '93: Proceedings of the*  
18 *1993 ACM SIGSOFT international symposium on Software testing and analysis* (New York,  
19 NY, USA, 1993), ACM Press, pp. 2–10.
- 20 [80] XIE, Q., AND MEMON, A. M. Automated model-based testing of community-driven open  
21 source GUI applications. In *ICSM '06: Proceedings of the 22nd IEEE International Con-*  
22 *ference on Software Maintenance* (Washington, DC, USA, 2006), IEEE Computer Society,  
23 pp. 145–154.

- 1 [81] XIE, Q., AND MEMON, A. M. Studying the characteristics of a ‘good’ GUI test suite. In  
2 *Proceedings of the 17th IEEE International Symposium on Software Reliability Engineering*  
3 *(ISSRE 2006)* (Nov. 2006), IEEE Computer Society Press.
- 4 [82] XIE, Q., AND MEMON, A. M. Designing and comparing automated test oracles for GUI-  
5 based software applications. *ACM Transactions on Software Engineering and Methodology*  
6 *16*, 1 (2007), 4.
- 7 [83] XIE, T., AND NOTKIN, D. Mutually enhancing test generation and specification inference.  
8 In Petrenko and Ulrich [59], pp. 60–69.
- 9 [84] XIE, T., AND NOTKIN, D. Tool-assisted unit-test generation and selection based on opera-  
10 tional abstractions. *Autom. Softw. Eng.* *13*, 3 (2006), 345–371.
- 11 [85] YILMAZ, C., COHEN, M. B., AND PORTER, A. Covering arrays for efficient fault character-  
12 ization in complex configuration spaces. In *Proceedings of the 2004 international symposium*  
13 *on Software testing and analysis* (2004), pp. 45–54.
- 14 [86] YUAN, X., AND MEMON, A. M. Using GUI run-time state as feedback to generate test  
15 cases. In *ICSE ’07: Proceedings of the 29th International Conference on Software Engineer-*  
16 *ing* (Minneapolis, MN, USA, May 23–25, 2007), IEEE Computer Society, pp. 396–405.