# Graphical user interface (GUI) testing: Systematic mapping and repository

**Published in:**
Information and Software Technology

**Document Version:**
Peer reviewed version

**Queen's University Belfast - Research Portal:**
Link to publication record in Queen's University Belfast Research Portal

# Graphical User Interface (GUI) Testing: Systematic Mapping and Repository

Ishan Banerjee[a], Bao Nguyen[a], Vahid Garousi[b], Atif Memon[a]

[a]*Department of Computer Science, University of Maryland, College Park, MD 20742, USA, {ishan, baonn, atif}@cs.umd.edu*
[b]*Electrical and Computer Engineering, University of Calgary, Calgary, Canada, vgarousi@ucalgary.ca*

**Abstract**

As the research area of "GUI testing" has matured, there has been an increase in the number of articles. More than 200 articles have appeared in this area since 1990. We study this body of knowledge using a systematic mapping (SM) in this paper. We define the term *GUI testing* as system testing of a software that has a graphical-user interface (GUI) front-end. Because system testing entails that the entire software system, including the user interface, be tested as a whole, during GUI testing, test cases, modeled as sequences of user input events, are created and executed on the software by exercising the GUI's widgets. As part of the SM, we pose three sets of research questions, define selection and exclusion criteria, and create a map of 136 articles. We share this map in a publicly accessible repository. We discuss future trends in GUI testing, and stress that articles in this area should clearly present certain attributes of their work to help conduct similar SMs in the future.

*Keywords:* systematic mapping, web application, testing, paper repository, bibliometrics

## Contents

## 1. Introduction

Whenever the number of *primary studies*—reported in *articles* (we use the term *article* to include research papers, book chapters, dissertations, theses, published experimental results, and published demonstrations of techniques)—in an area grows very large, it is useful to summarize the body of knowledge and to provide an overview using a secondary study [81]. A secondary study [3, 4, 19, 52] aggregates and objectively synthesizes the outcomes of the primary studies. Because the synthesis needs to have some common basis for extracting attributes in the articles, a side effect of the secondary study is that it encourages researchers conducting and reporting primary studies to

improve their reporting standard of such attributes, which may include metrics, tools, study subjects, limitations, etc. Moreover, by "mapping the research landscape," a secondary study helps to identify sub-areas that need more primary studies.

In the field of Software Engineering (SE), a systematic mapping (SM) study is a well-accepted method to identify and categorize research literature [20, 81]. An SM [3, 12, 28, 35, 56, 79, 82] study focuses on building classifications schemes and the results show frequencies of articles for classes within the scheme. These results become one of the outputs of the SM in the form of a database or *map* that can be a useful descriptive tool itself. An SM uses established searching protocols and has rigorous inclusion/exclusion criteria.

In this paper, we leverage the guidelines set by Petersen *et al.* [81] and Kitchenham *et al.* [55] to create an SM for the area of *GUI testing*. We define the term GUI testing to mean that a GUI-based application, *i.e.*, one that has a graphical-user interface (GUI) front-end, is tested solely by performing sequences of events (*e.g.*, *"click on button"*, *"enter text"*, *"open menu"*) on GUI *widgets* (*e.g.*, *"button"*, *"text-field"*, *"pull-down menu"*). In all but the most trivial GUI-based systems, the space of all possible event sequences that may be executed is extremely large, in principle infinite (consider the fact that a user of MS Word can click on the *File* menu an unlimited number of times). All GUI testing techniques are in some sense *sampling* the input space, either manually [9, 78] or automatically [70, 95]. In the same vein, techniques that develop a GUI *test oracle* [11]—a mechanism that determines whether a GUI executed correctly for a test input—are based on sampling the output space; examining the entire output, pixel by pixel, is simply not practical [89, 98]. Techniques for evaluating the adequacy of GUI test cases provide some metrics to quantify the test cases [71, 103, 108]. And techniques for re-

gression testing focus on retesting the GUI software after modifications [49, 66, 96].

The above is just one possible classification of GUI testing techniques. The goal of our SM is to provide a much more comprehensive classification of the over 200 articles that have appeared in the area since 1990. Given that now there are regular events such as the International Workshop on TESTing Techniques & Experimentation Benchmarks for Event-Driven Software (TESTBEDS) [93] in the area, we expect this number to increase. We feel that this is an appropriate time to discuss trends in these articles and provide a synthesis of what researchers think are limitations of existing techniques and future directions in the area. We also want to encourage researchers who publish results of primary studies to improve their reporting standards, and include certain attributes in their articles to help conduct secondary studies. Considering that many computer users today use GUIs exclusively and have encountered GUI-related failures, research on GUIs and GUI testing is timely and relevant.

There have already been 2 smaller, preliminary secondary studies on GUI testing. Hellmann *et al.* [90] presented a literature review of test-driven development of user interfaces; it was based on a sample of 6 articles. Memon *et al.* [67] presented a classification of 33 articles on model-based GUI test-case generation techniques. To the best of our knowledge, there are no other secondary studies in the area of GUI testing.

In our SM, we study a total of 213 articles. We formulate 3 sets of research questions pertaining to the research space of GUI testing, demographics of the studies and authors, and synthesis and interpretation of findings. We describe the mechanisms that we used to locate the articles and the set of criteria that we applied to exclude a number of articles; in all we classify 136 articles. Our most important findings suggest that there is an increase in the number of articles in the area; there has been lack of evaluation and validation, although this trend is changing; there is insufficient focus on mobile platforms; new techniques continue to be developed and evaluated; evaluation subjects are usually non trivial, mostly written in Java, and are often tested using automated model-based tools; and by far a large portion of the articles are from the US, followed by China.

We have published our SM as an online repository on Google Docs [94]. Our intention is to periodically update this repository, adding new GUI testing articles as and when they are published. In the future, we intend to allow authors of articles to update the repository so that it can become a "live" shared resource maintained by the wider GUI testing community.

The remainder of this paper is structured as follows. Section 2 presents background and related work. Section 3 presents our goals and poses research questions. The approach that we used to select articles is presented in Section 4. Section 5 presents the process used for constructing the systematic map. Sections 6, 7, and 8 present the results of the systematic mapping. Finally, Section 9 concludes with remarks and discussion.

## 2. Background and Related Work

In this section, we present more details of GUI testing. We also summarize the 14 secondary studies that have been reported in the broader area of software testing. Finally, because we are sharing the data artifacts produced from our SM in an online repository, we discuss others' efforts to do the same.

**GUI Testing:** As computers play an increasingly important role aiding end-users, researchers, and businesses in today's internetworked world, the class of software that has a graphical user interface (GUI) front-end has become ubiquitous [72, 39, 87]. A GUI takes events (mouse clicks, selections, typing in textfields) as input from users, and then changes the state of its widgets. GUIs have become popular because of the advantages this "event-handler architecture" offers to both developers and users [33, 105]. From the developer's point of view, the event handlers may be created and maintained fairly independently; hence, complex system may be built using these loosely coupled pieces of code. From the user's point of view, GUIs offer many degrees of usage freedom, i.e., users may choose to perform a given task by inputting GUI events in many different ways in terms of their type, number and execution order.

Testing and Quality Assurance (QA) is becoming increasingly important for GUIs as their functional correctness may affect the quality of the entire system in which the GUI operates. Software testing is a popular QA technique employed during software development and deployment to help improve its quality [43, 63]. During software testing, test cases are created and executed on the software. One way to test a GUI is to execute each event individually and observe its outcome, thereby testing each event handler in isolation [70]. However, the execution outcome of an event handler may depend on its internal state, the state of other entities (objects, event handlers) and the external environment. Its execution may lead to a change in its own state or that of other entities. Moreover, the outcome of an event's execution may vary based on the sequence of preceding events seen thus far. Consequently, in GUI testing, each event needs to be tested in different states. GUI testing therefore involves generating and executing sequences of events [104, 105]. Most of the articles on test generation that we classify in our SM consider the event-driven nature of GUI test cases, although few mention it explicitly.

**Secondary studies in software testing:** There have been 14 reported secondary studies in different areas of software testing, 2 related to GUI testing. We list these studies in Table 1 along with some of their attributes. For example, the "number of articles" column (No.) shows that the number of primary studies analyzed in each study varied from 6 (in [90]) to 264 (in [52]), giving some idea of the comprehensiveness of the studies.

Of particular interest to us are the SMs and structured literature reviews (SLRs). An SLR analyzes primary studies, reviews them in depth and describes their methodology and results. SLRs are typically of greater depth than SMs. Often, SLRs include an SM as a part of the study. Typically SMs and SLRs formally describe their search protocol and inclusion/exclusion criteria. We note that SMs and SLRs have re-

Table 1: 14 Secondary Studies in Software Testing

| Type | Secondary Study Area | No. | Year | Ref. |
|---|---|---|---|---|
| SM | Non-functional search-based soft. testing | 35 | 2008 | [3] |
| | SOA testing | 33 | 2011 | [79] |
| | Requirements specification and testing | 35 | 2011 | [12] |
| | Product lines testing | 45 | 2011 | [28] |
| SLR | Search-based non-functional testing | 35 | 2009 | [4] |
| | Search-based test-case generation | 68 | 2010 | [5] |
| Survey | Object oriented testing | 140 | 1996 | [19] |
| | Testing techniques experiments | 36 | 2004 | [53] |
| | Search-based test data generation | 73 | 2004 | [65] |
| | Combinatorial testing | 30 | 2005 | [42] |
| | Symbolic execution for software testing | 70 | 2009 | [83] |
| Taxonomy | Model-based GUI testing | 33 | 2010 | [67] |
| Lit rev. | Test-driven development of user interfaces | 6 | 2010 | [90] |
| Analysis/survey | Mutation testing | 264 | 2011 | [52] |

cently started appearing in the area of software testing. There are four SMs: product lines testing [28], SOA testing [79], requirements specification and testing [12] and non-functional search-based software testing [3]. There are two SLRs – search-based non-functional testing [4] and search-based test-case generation [5].

The remaining 8 studies are "surveys", "taxonomies", "literature reviews", and "analysis and survey", terms used by the authors themselves to describe their studies.

**Online Article Repositories in SE:** Authors of a few recent secondary studies have developed online repositories to supplement the study. This is a large undertaking as even after the study is published, these repositories are updated regularly, typically every 6 months to a year. Maintaining and sharing such repositories provides many benefits to the broader community. For example, they are valuable resources for new researchers (e.g., PhD students) and for other researchers aiming to do additional secondary studies.

For example, Mark Harman and his team have developed and shared two online repositories, one in the area of mutation testing [52], and another in the area of search-based software engineering (SBSE) [62, 92]. The latter repository is quite comprehensive and has 1014 articles as of Mar. 2012, a large portion of which are in search-based testing.

## 3. Goals, Questions, and Metrics

We use the Goal-Question-Metric (GQM) paradigm [13] to form the goals of this SM, raise meaningful research questions, and carefully identify the metrics that we collect from our data and how we use them to create our maps. The goals of this study are:

**G1:** To classify the nature of articles in the area of GUI testing, whether new techniques are being developed, whether they are supported by tools, their weaknesses and strengths, and to highlight and summarize the challenges and lessons learned.

**G2:** To understand the various aspects of GUI testing (e.g., test creation, test coverage) that are being researched.

**G3:** To study the nature of evaluation, if any, that is being conducted, the tools being used, and subject applications.

**G4:** To identify the most active researchers in this area and their affiliations, and identify the most influential articles in the area.

**G5:** To determine the recent trends and future research directions in this area.

Goals **G1**, **G2**, and **G3** are all related to understanding the trends in GUI testing *research and evaluation* being reported in articles. These goals lead to our first set of research questions. Note that as part of the research questions, we include the metrics (underlined) that we collect for the SM.

**RQ 1.1:** *What types of articles have appeared in the area?* For example, we expect some articles that present new techniques, others that evaluate and compare existing techniques.

**RQ 1.2:** *What test data generation approaches have been proposed?* For example, some test data may be obtained using manual approaches, other via automated approaches.

**RQ 1.3:** *What type of test oracles have been used?* A test oracle is a mechanism that determines whether a test case passed or failed for a given test input. A test case that does not have a test oracle is of little value as it will never fail. We expect some test cases to use a manual test oracle, *i.e.*, manual examination of the test output to determine its pass/fail status. Other test cases may use an automated test oracle, in which the comparison between expected and actual outputs is done automatically.

**RQ 1.4:** *What tools have been used/developed?* We expect that some techniques would have resulted in tools; some are based on existing tools. Here we want to identify the tools and some of their attributes, *e.g.*, execution platform.

**RQ 1.5:** *What types of systems under test (SUT) have been used?* Most new techniques need to be evaluated using some software subjects or SUTs. We want to identify these SUTs, and characterize their attributes, *e.g.*, platform (such as mobile, web), size in lines of code (LOC).

**RQ 1.6:** *What types of evaluation methods have been used?* We expect that some techniques would have been evaluated using the type and amount of code that they cover, others using the number of test cases they yield, and natural or seeded faults they detected.

**RQ 1.7:** *Is the evaluation mechanism automated or manual?* A new technique that can be evaluated using automated mechanisms (*e.g.*, code coverage using code instrumentation) makes it easier to replicate experiments and conduct comparative studies. Widespread use of automatic mechanisms thus allows the research area to encourage experimentation.

To answer all the above questions, we carefully examine the articles, collect the relevant metrics, create classifications replying explicitly on the data and findings reported in the articles, and obtain frequencies when needed. All the metrics are objective, *i.e.*, we do not offer any subjective opinions to answer any of these questions.

Goals **G4** and parts of **G1** and **G5** are concerned with understanding the *demographics and bibliometrics* of the articles and

authors. These goals lead to our second set of research questions.

**RQ 2.1:** *What is the annual articles count?*

**RQ 2.2:** *What is the article count by venue type?* We expect the most popular venues to be conferences, workshops, and journals.

**RQ 2.3:** *What is the citation count by venue type?*

**RQ 2.4:** *What are the most influential articles in terms of citation count?*

**RQ 2.5:** *What were the venues with highest articles count?*

**RQ 2.6:** *What were the venues with highest citation count?*

**RQ 2.7:** *Who are the authors with the largest number of articles?*

**RQ 2.8:** *What are the author affiliations, i.e., do they belong to academia or industry?*

**RQ 2.9:** *Which countries have produced the most articles?*

Again, we observe that the above questions may be answered by collecting objective metrics from the articles.

Goals **G5** and parts of **G1** are concerned with the *recent trends, limitations, and future research directions* in the area of GUI testing; we attain these goals by studying recent articles, the weaknesses/strengths of the reported techniques, lessons learned, and future trends. More specifically, we pose our third set of research questions.

**RQ 3.1:** *What limitations have been reported?* For example, some techniques may not scale for large GUIs.

**RQ 3.2:** *What lessons learned are reported?*

**RQ 3.3:** *What are the trends in the area?* For example, new technologies may have prompted researchers to focus on developing techniques to meet the needs of the technologies.

**RQ 3.4:** *What future research directions are being suggested?*

Due to the nature of the questions, their answers may be based on opinions of the original authors who conducted the primary studies.
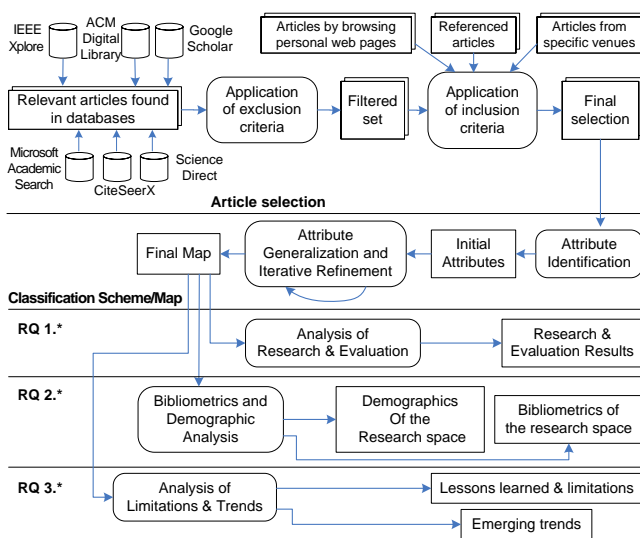


Figure 1: Protocol Process Guiding this SM.

Having identified the goals for this work, linking them to research questions, and identifying the metrics that we collect, we have set the stage for the SM. The remainder of this paper is based on the protocol that lies at the basis of this SM; it is outlined in Figure 1. Note that the protocol distinguishes five phases that are described in Sections 4–8. More specifically, we describe the process of article selection in Section 4, map construction in Section 5, and address research questions **RQ 1.*** in Section 6, **RQ 2.*** in Section 7, and **RQ 3.*** in Section 8.

## 4. Article Selection

As can be imagined, article selection is a critical step in any secondary study. Indeed, it lays the foundation for the synthesis of all of its results. Consequently, in any secondary study, article selection must be explained carefully so that the intended audience can interpret the results of the study keeping in mind the article selection process. In this work, the articles were selected using a three step process using guidelines presented in previous systematic mapping articles [107, 55, 81]: (1) *article identification*, done using digital libraries and search engines, (2) *definition and application of exclusion criteria*, which exclude articles that lie outside the scope of this study, and (3) *definition and application of inclusion criteria*, which target specific resources and venues that may have been missed by the digital libraries and search engines to hand-pick relevant articles. These steps are illustrated in the top part of Figure 1. We now expand upon each step.

**Step 1: Article Identification:** We started the process by conducting a keyword-based search to extract a list of articles from the following digital libraries and search engines: IEEE Xplore,[1] ACM Digital Library,[2] Google Scholar,[3] Microsoft Academic Search,[4] Science Direct,[5] and CiteSeerX[6]. The following keywords were used for searching: *GUI testing*, *graphical user interface testing*, *UI testing*, and *user interface testing*; we looked for these keywords in article titles and abstracts. This step yielded 198 articles forming the initial pool of articles.

**Step 2: Exclusion Criteria:** In the second step of the process, the following set of exclusion criteria were defined to exclude articles from the above initial pool. **C1:** languages other than English, **C2:** not relevant to the topic, and **C3:** that did not appear in the published proceedings of a conference, symposium, or workshop, or did not appear in a journal or magazine.

These criteria were then applied by defining application procedures. It was fairly easy to apply criterion **C1** and **C3**. For criterion **C2**, a voting mechanism was used amongst us (the authors) to assess the relevance of articles to GUI testing. We focused on the inclusion of articles on functional GUI testing; and excluded articles on non-functional aspects of GUIs, such

---

[1] http://ieeexplore.ieee.org/
[2] http://dl.acm.org/
[3] http://scholar.google.com/
[4] http://academic.research.microsoft.com/
[5] http://www.sciencedirect.com/
[6] http://citeseer.ist.psu.edu

as stress testing GUI applications [2] and GUI usability testing [75]. Application of the above exclusion criteria resulted in a filtered set of 107 articles.

**Step 3: Inclusion Criteria:** Because search engines may miss articles that may be relevant to our study, we supplemented our article set by manually examining the following three sources: (1) web pages of active researchers, (2) bibliography sections of articles in our filtered pool, and (3) specific venues.

These sources led to the definition of 3 corresponding inclusion criteria. Application of the first two criteria was straightforward. For the third criterion, the specific venue that had not been indexed by the popular search engines was TESTing Techniques & Experimentation Benchmarks for Event-Driven Software (TESTBEDS), which is a relatively new workshop. Application of the 3 inclusion criteria resulted in the final pool of articles containing 136 articles.

**Our Final Article Set:** Figure 2 shows the distribution of the 230 articles analyzed during this study. The dark shaded part of each horizontal bar shows the number that we finally included, forming a total of 136 articles. A few articles are classified as "Unknown" because, despite numerous attempts, we were unable to obtain them. In summary, we have included all articles presented at all venues that print their proceedings or make them available digitally.
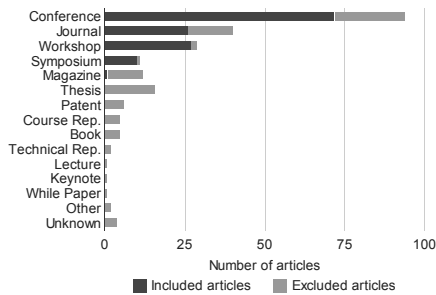


Figure 2: Total Articles Studied = 230; Final Included = 136

## 5. Map Construction

As mentioned earlier, a map is the tool used for classification of the selected articles. Construction of the map is a complex and time-consuming process. Indeed the map that we have made available in a publicly accessible repository is one of the most important contributions of our work. Fortunately, because we use the GQM approach, we already have research questions and metrics; we use the metrics as a guide for map construction. For RQ 1.*, we need to collect the metrics: "types of articles," "test data generation approaches," "type of test oracles," "tools," "types of SUT," "types of evaluation methods," and "evaluation mechanism." This list in fact forms a set of classes of *attributes* of the articles. We define these attributes in this section and present the map structure; with this map (also called attribute framework [85]), the articles under study can be characterized in a comprehensive fashion.

The map was created in an iterative manner. In the first iteration, all articles were analyzed and terms which appeared to be of interest or relevance for a particular aspect (e.g., 'subject under test', 'testing tool'), were itemized. This itemization task was performed by all of us. To reduce individual bias, we did not assume any prior knowledge of any attributes or keywords. The result after analyzing all articles was a large set of initial attributes. After the initial attributes were identified, they were generalized. This was achieved through a series of meetings. For example, under "test data generation approaches," the attributes 'finite-state machine (FSM)-based' method and 'UML-based' method were generalized to 'model-based' method.

Defining attributes for "types of articles" was quite complex. As one can imagine, there are innumerable ways of understanding the value of a research article. To make this understanding methodical, we defined two *facets*—specific ways of observing a subject—which helped us to systematically understand the contribution and research value of each article. The specific facets that we used, *i.e.*, *contribution* and *research* were motivated from [81].

The resulting attributes for each facet were documented, yielding a map that lists the aspects, attributes within each aspect, and brief descriptions of each attribute. This map forms the basis for answering the research questions **RQ 1.***.

Similarly, for RQ 2.* we need the following metrics: "annual articles count," "article count by venue type," "citation count by venue type," "citation count," "citation count by venue," "venues with highest article counts," "authors with maximum articles," "author affiliations," and "countries." The first two metrics were obtained directly from our spreadsheet. The remaining metrics lead us to develop our second map. As before, the map lists the attributes and brief descriptions of each attribute. This map forms the basis for answering the research questions **RQ 2.***.

Finally, for RQ 3.*, we need to collect the metrics: "limitations," "lessons learned," "trends," and "future research directions." This led us to develop our third map, which forms the basis for answering the research questions **RQ 3.***. The final map used in this research for all questions is shown in Figure 3.

## 6. Mapping Research & Evaluation

We are now ready to start addressing our original research questions **RQ 1.1** through **RQ 1.7**.

**RQ 1.1:** *What types of articles have appeared in the area?* As discussed earlier in Section 5, we address this question using two facets, primarily taken from [81]. The *contribution facet*—test method, test tool, test model, metric, process, challenge, empirical study—broadly categorizes the type of the article. On the other hand, the *research facet*—solution proposal, validation, evaluation research, experience, philosophical and opinion articles—broadly categorizes the nature of research work presented in the article. It helps understand the nature of research exploration done in the article. Every article has been attributed at least one category. Some articles have been placed in more than one category. For example, Belli [14] presents a

| | | | |
|---|---|---|---|
| **RQ 1.1: Type of article** | *Contribution Facet* | Test method/technique (B) | Article describes new technique or improves upon an existing one |
| | | Test tool (B) | Article focuses on testing tool and evaluates its applicability. |
| | | Test model (B) | Article introduces new modeling technique or is based on use of model |
| | | Metric (B) | Article describes new metric for evaluating testing techniques |
| | | Process (B) | Article describes software testing process or life-cycle |
| | | Challenge (B) | Article discusses challenges in certain areas of GUI testing |
| | | Empirical study (B) | Article is an empirical study of technique |
| | *Research Facet* | Solution proposal (B) | New solution; applicability shown via example or line of argument |
| | | Validation research (B) | Novel technique demonstrated in lab with experiment |
| | | Evaluation research (B) | Comprehensive experimental evaluation of technique |
| | | Experience article (B) | Personal experience with GUI testing of authors |
| | | Philosophic article (B) | Sketch new way of looking at existing things. |
| | | Opinion article (B) | Opinion of authors on the goodness of techniques. |
| **RQ 1.2** | Test data generation | Capture/replay (B) | Capture/replay was used to generate test cases |
| | | Model based (B) | GUI model was used to generate test cases |
| | | Model name (S) | Name of model used (if model-based) |
| | | Random testing (B) | Test cases were generated randomly |
| **RQ 1.3** | Test oracle | State reference (B) | GUI state information was used as oracle |
| | | Crash testing (B) | SUT crash was used to identify faults |
| | | Formal specification (B) | Formal specification of SUT was used as oracle |
| | | Manual verification (B) | Result of test case execution was manually verified |
| | | Multiple oracles (B) | More than one oracle was used in same test run |
| **RQ 1.4** | Testing tools | Tool proposed (S) | Name of new tool introduced in an article |
| | | Tool used (S) | Name of existing or third party tool used in an article |
| | | Programming language (S) | Programming language used in developing the tool |
| **RQ 1.5** | System under test | Number of SUT(s) (N) | Number of SUT(s) used in the article |
| | | Size (LOC) (N) | Number of lines of code in SUT |
| | | Programming language (S) | Programming language of the SUT |
| | | GUI technology (S) | GUI SDK or library used in SUT |
| | | Small/large scale (S) | Qualitative assesment of the size of SUT |
| **RQ 1.7** | Evaluation Automation | Automated (B) | Automated test case execution was used in article |
| | | Manual (B) | Manual test case execution was used in article |
| | | None (B) | Test cases were not executed |
| **RQ 2.*** | Demographic Information | Authors (S) | Name of all contributing authors |
| | | Authors' country (S) | Country from which author published the article |
| | | Authors' affiliations (E) | Are the authors from academia, industry or a mix of both |
| | | Venue (S) | Where it was published |
| | | Year (N) | Year of publication |
| | | Citation count (N) | Number of times this work has been cited, per year, as of July 2012 |
| **RQ 3.*** | Limitations & Future | Limitations (S) | Limitations noted by article authors |
| | | Lessons learned (S) | Lessons learned |
| | | Future research (S) | Future research directions |

* B = Boolean, E = Enumerated, N = Numeric, S = String

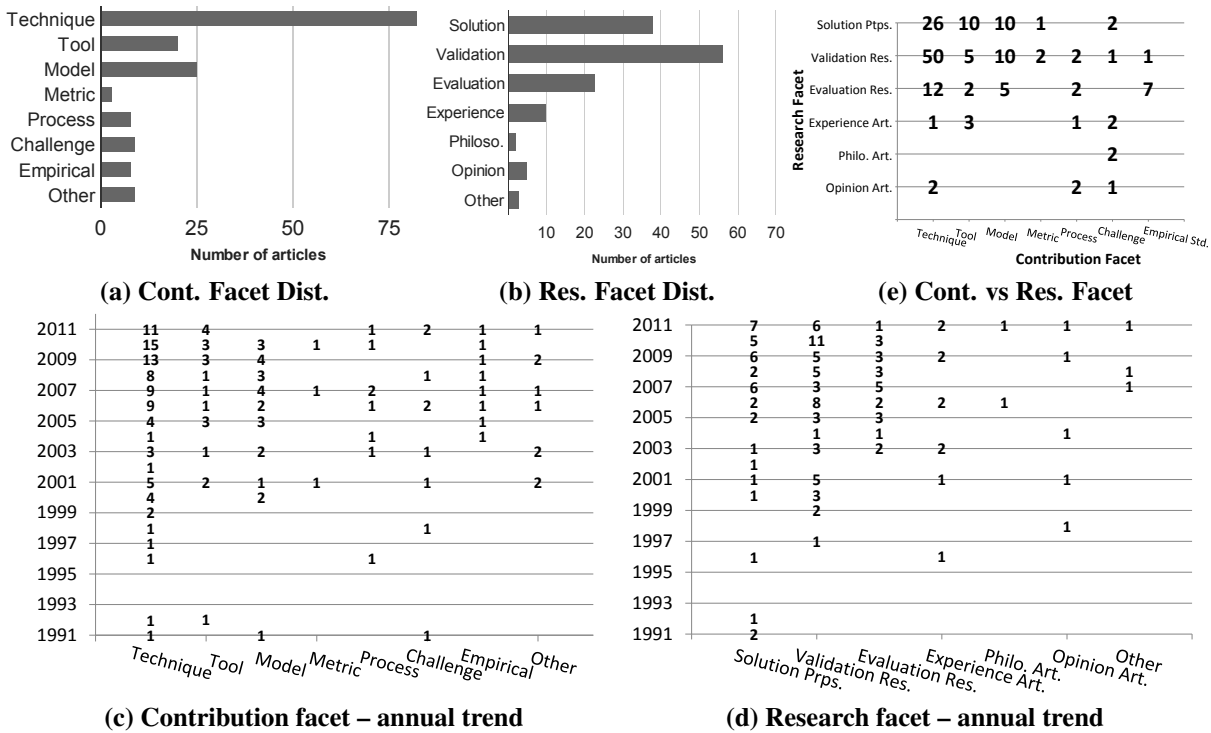Figure 3: The Final Map Produced by and Used in this Research

**(a) Cont. Facet Dist.**

**(b) Res. Facet Dist.**

**(e) Cont. vs Res. Facet**

**(c) Contribution facet – annual trend**

**(d) Research facet – annual trend**

Figure 4: Data for **RQ 1.1**.

testing technique based on FSMs. This article is placed under both 'test method' as well as 'test model' in contribution facet.

Figure 4(a) shows the contribution facet for all the 136 articles. The y-axis enumerates the categories, and the x-axis shows the number of articles in each category. Most articles (90 articles) have contributed towards the development of new or improved testing techniques. Few articles have explored GUI testing metrics, or developed testing processes. Figure 4(c) shows an annual distribution of the contribution facet. The y-axis enumerates the period 1991-2011, the x-axis enumerates the categories, the integer indicates the number of articles in each category for a year. During the period 1991-2000, most of the work focused on testing techniques. On the other hand, during 2001-2011, articles have contributed to various categories. This trend is likely owing to the rising interest in GUI testing in the research community.

Figure 4(b) shows the research facet for all the 136 articles. Most articles propose solutions, conduct various types of experiments to validate techniques. There are very few philosophical or opinion articles. Figure 4(d) shows an annual distribution of the research facet. From the figure, there is an increasing number of articles in recent years, with most articles in solution proposal, validation and evaluation research. In the year 2011, the largest number of articles were on validation research, a promising development, showing that researchers are not only proposing novel techniques, but they are also supporting them with lab experiments.

To get a better understanding of the contribution and research focus of each article, we also visualize the relationship between the research and contribution facets in Figure 4(e) The y-axis

enumerates the research facet categories; the x-axis enumerates the contribution facet categories. The intersection of each pair of categories is an integer whose value corresponds to the number of articles at that point. Work on exploring new and improved techniques dominate with focus on validation research with 46 articles. A small but noticeable amount of work has been done on empirical research focusing on extensive evaluation of techniques with 7 articles.

**RQ 1.2:** *What test data generation approaches have been proposed?* Of the 136 articles, 123 articles reported generation of test artifacts. For example, Ariss *et al.* [9] uses capture/replay and model-based methods for testing Java applications. Figure 5(a) shows the distribution of test data generation methods. The x-axis shows the number of articles for each method; y-axis enumerates the methods. Model-based (72 articles) and capture/replay based (23 articles) methods are most common. The remaining 37 articles use less popular methods such as symbolic execution [36], formal method [91], AI planning [70], statistical analysis [88], etc.

Since model-based methods are commonly used, Figure 5(b) shows the composition of these 72 articles. The x-axis shows the number of articles using a model, y-axis enumerates the models. Models such as event flow graph (EFG), finite state machine (FSM) were most common. There are 25 articles which use less common models such as Probabilistic model [15], Function trees [61], etc.

**RQ 1.3:** *What type of test oracles have been used?* We remind the reader that a test oracle is a mechanism that determines if a test case passed or failed. As Figure 6 shows, state reference
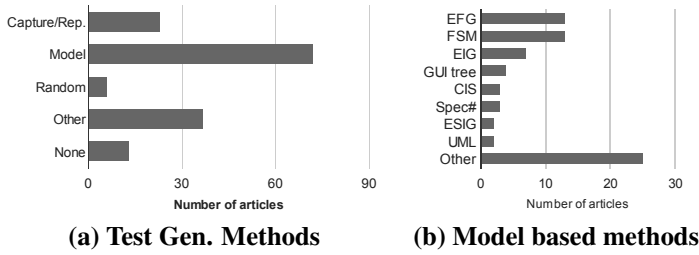
7

**(a) Test Gen. Methods**　　　**(b) Model based methods**
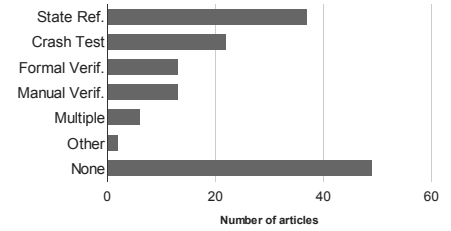
Figure 5: Data for **RQ 1.2**.



Figure 6: Data for **RQ 1.3**: Oracle Type

(37 articles) is the commonly used oracle. In this method the state of the GUI is extracted while the SUT is executing, and is stored. At a later time, this state may be compared with another execution instance for verification [98, 89]. SUT crash testing is another popular oracle (22 articles). In this method, if the SUT crashes during the execution of a test case, then the test case is marked as failed. The 'crash' state of the SUT is thus an oracle [103, 7]. Formal verification (13 articles) methods use a model or specification to verify the correctness of the output of a test case [69, 91]. Manual verification (13 articles) is also used. In this method a human tester is involved in verifying the result of executing a test case [95, 58].

We observed that a large number of articles (49 articles) did not use a test oracle. Of these, 13 articles are experience, philosophical or opinion articles and do not require a test oracle for evaluation. The remaininga 36 articles are solution proposal, validation or evaluation but do not use a test oracle (*e.g.*, [36, 49]).

**RQ 1.4:** *What tools have been used/developed?* Testing of GUI based applications typically require the use of tools. A *tool*, for the purpose of this paper, is a set of well defined, packaged, distributable software artifacts which is used by a researcher to evaluate or demonstrate a technique. Test scripts, algorithm implementations and other software components used for conducting experiments, which were not named or did not appear to be easily distributable, have not been considered as tools.

A tool is considered as a *new tool* if it has been developed specifically for use in an article. A tool is considered as an *existing tool* if it has been developed by the author in a previous work or has been developed by a third party – commercially available, open source, etc.

Figure 7(a) shows the composition of new and existing tools used for all 136 articles. It can be seen that 32 articles (23.52%) introduced a new tool only, 48 articles (35.29%) used an existing tool only, 29 articles (21.32%) used both new and existing tools, whereas 27 articles (19.85%) did not use a clearly defined tool. From this figure it can be seen that most articles (109 articles) used one or more tools. Certain articles, such as experience, philosophical and opinion articles, for example by Robinson *et al.* [84], did not require a tool.

From the 109 articles that used a tool, a total of 112 tools were identified. Note that a tool may have been used in more than one article. Similarly, an article may have used more than one tool. Figure 7(b) shows the ten most popular tools and their usage count. The x-axis shows the number of articles where

the tool was used, y-axis enumerates the 10 most popular tools. GUITAR [1], which ranks highest, has been used in 22 articles. 91 tools were used in only 1 article, 15 tools were used in 2 articles and so forth. Only 1 tool, GUITAR [1], was used in 22 articles.

New GUI testing tools were described in 61 articles. Figure 7(c) shows the distribution of programming languages in which the tools were developed. The x-axis shows the number of articles in which a new tool was developed in a particular language, y-axis enumerates the languages. From the figure, Java is by far the most popular choice with 23 articles.

**RQ 1.5:** *What types of systems under test (SUT) have been used?* Of the 136 articles, 118 reported the use of one or more SUT. Note that an SUT may have been used in different articles, conversely, more than one SUT may have been used in an article. Figure 8(a) shows the number of SUTs that were used in each article. This figure helps us understand how many SUTs are typically used by researchers to evaluate their techniques. The x-axis enumerates the SUT count from 1–6 and ≥ 7. The y-axis shows the number of articles using a given number of SUTs. From the figure, it can be seen that out of 136 articles, 118 used one or more SUTs. Only 1 SUT was used in 64 articles, while a small number of articles (5) [31, 32, 47, 48, 61] used 7 or more SUTs. A total of 18 articles (e.g., [45]) did not use any SUT.

Figure 8(b) shows the programming language of SUTs reported by 71 articles. The x-axis enumerates the common languages, y-axis shows the number of articles for each language. We see that Java applications are by far the most common SUT with 48 articles using Java based SUT(s) [105, 34, 46]. C/C++ [84, 26] and .NET [25, 6] based SUTs have been used in 16 articles. The remaining 7 SUTs are based on MATLAB [29], Visual Basic [59] and Objective C [22].

SUTs also differed in their underlying development technology. For example, some SUTs were developed using Abstract Window Toolkit (AWT), while others were developed using Swing. Classifying their development technology helps understand the experimental environment that are the focus of new GUI testing techniques. We found that Java Swing is by far the most common technology, with 25 out of 36 articles using an SUT based on Swing. This is consistent with a large number of SUTs being based on Java.

SUTs reported by researchers varied in size, in terms of lines-of-code (LOC) – some less than 1,000 lines while some more than 100,000 lines. Only 28 articles out of 136 reported this
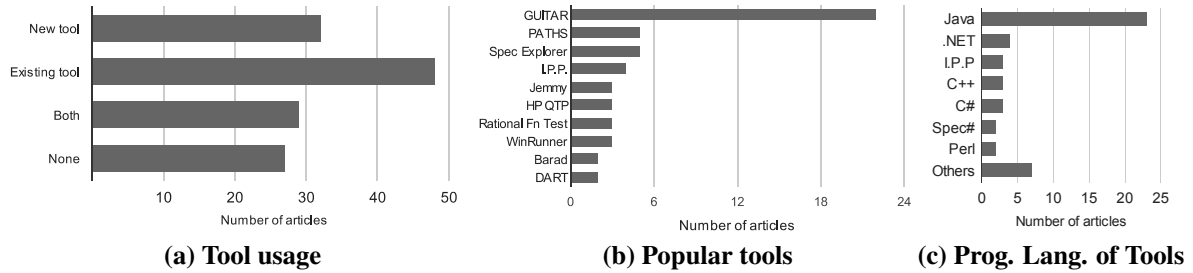
8

**(a) Tool usage**  **(b) Popular tools**  **(c) Prog. Lang. of Tools**

Figure 7: Data for **RQ 1.4**.



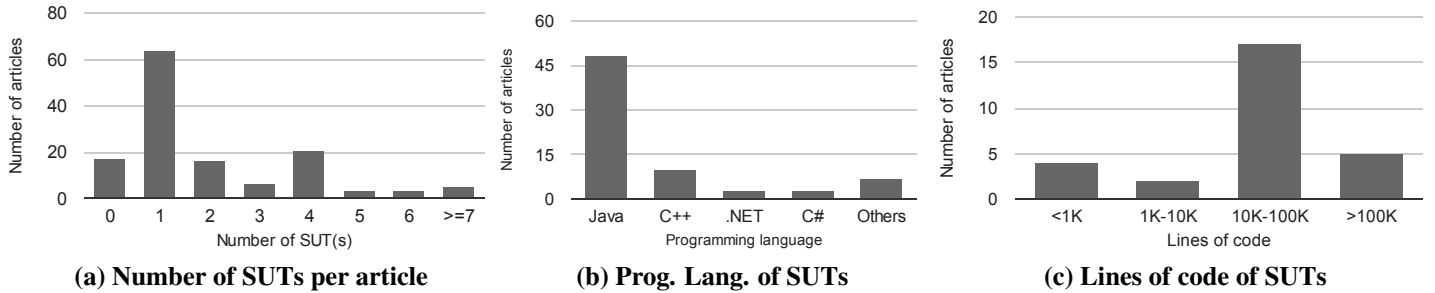**(a) Number of SUTs per article**  **(b) Prog. Lang. of SUTs**  **(c) Lines of code of SUTs**

Figure 8: Data for **RQ 1.5**.

information. Figure 8(c) shows the cumulative LOC of SUTs used in each article. The x-axis enumerates ranges of LOC, y-axis shows the number of articles in each range. Most articles used SUTs in the range 10,000-100,000 (17 articles). Only 5 articles [10, 25, 44, 66, 106] used SUTs with LOC totaling more than 100,000 lines.

The SUTs were also classified as *large-scale* or *small-scale*. This classification helps us understand if some articles used small or *toy* SUTs. SUTs such as commercially available software – Microsoft WordPad [68], physical hardware such as vending machines [51], mobile phones [54] and open source systems [66] have been classified as large-scale systems. SUTs such as a set of GUI windows [73], a set of web pages [45], small applications developed specifically for demonstration [24, 37] have been classified as small-scale SUTs. Of the 118 articles which used one or more SUTs, 89 articles (75.42%) used a large-scale SUT.

**RQ 1.6:** *What types of evaluation methods have been used?* Many articles studied in this SM focused on the development of new GUI testing techniques. The techniques developed in the article were evaluated by executing test cases on an SUT. Different methods and metrics were applied to determine the effectiveness of the testing technique.

A total of 119 articles reported one or more evaluation methods. Figure 9(a) shows the distribution of evaluation methods. The x-axis shows the count of articles in each method, eleven evaluation methods are enumerated on the y-axis. For example, 47 articles demonstrated the feasibility of the technique using a simple example. Figure 9(b) shows the metrics used in the evaluation. The x-axis shows the number of articles for each evaluation metric, y-axis enumerates evaluation metrics. Out of 136 articles, 75 articles specified an evaluation metric. Of

these, the *number of faults detected* was the common metric (32 articles).

The number of generated test cases were reported and used in 52 of the 136 articles for the evaluation process. Figure 9(c) shows the number of test cases used. The x-axis enumerates ranges of test case counts, y-axis shows the number of articles in each range. Most articles used less than 1,000 test cases. Four articles used more than 100,000 test cases [21, 97, 99, 105].

**RQ 1.7:** *Is the evaluation mechanism automated or manual?* Of the 136 articles, 86 articles reported execution of test cases for evaluation, of which 72 reported automated test case execution, 11 articles reported manual test case execution while 3 articles [40, 64, 80] reported both automated and manual test case execution.

## 7. Mapping Demographics

We now address the **RQ 2.*** research questions set, which is concerned with understanding the demographics of the articles and authors.

**RQ 2.1:** *What is the annual articles count?* The number of articles published each year were counted. The trend of publication from 1991 to 2011 is shown in Figure 10. An increasing trend in publication over years is observed. The two earliest articles in the pool were published in 1991 by Yip *et al.* [101, 102]. The total article counts of GUI testing articles in 2010 and 2011 was 19.

**RQ 2.2:** *What is the articles count by venue?* We classify the articles by venue type – conference, journal, workshop, symposium or magazine. Figure 11 shows that the number of con-
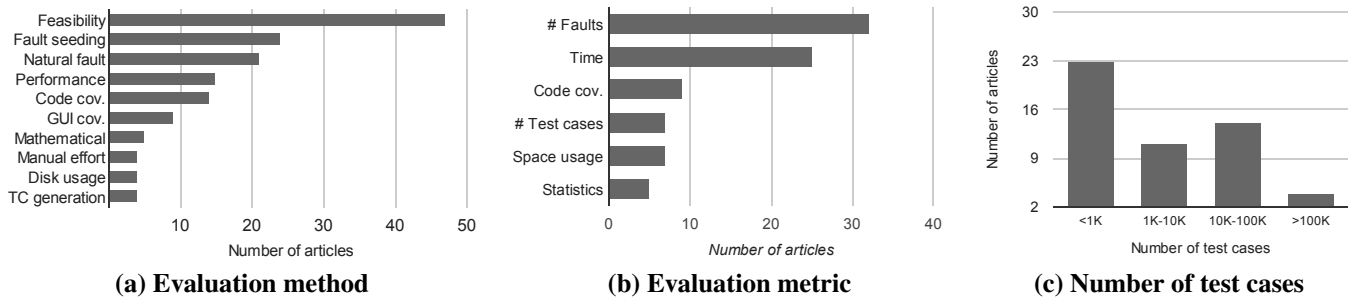
**(a) Evaluation method**



**(b) Evaluation metric**



**(c) Number of test cases**

Figure 9: Data for **RQ 1.6**.
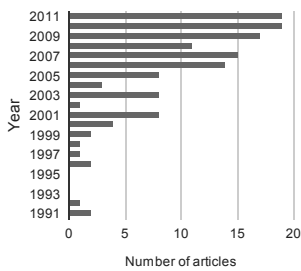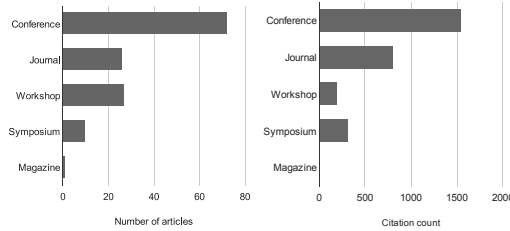


Figure 10: **RQ 2.1**: Annual Counts



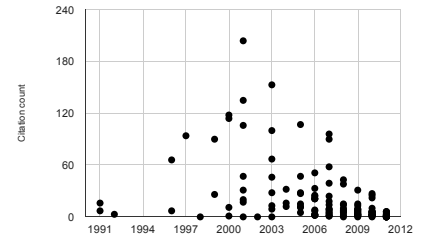Figure 11: **RQ 2.2** and **2.3**: Venue Types



Figure 12: **RQ 2.4**: Citations vs. Year

ference articles (72) are more than the articles in the other four categories combined (64).

**RQ 2.3:** *What is the citation count by venue type?* The number of citations for each article was extracted and aggregated for each venue type. Figure 11 shows the number of citations from different venue types. Conferences articles have received the highest citations at 1544.

**RQ 2.4:** *What are the most influential articles in terms of citation count?* This research question analyzes the relationship between the citations for each article and its year of publication. Figure 12 shows this data. The x-axis is the year of publication, and the y-axis is the number of citations. Each point in the figure represents an article.

The points for the recent articles (from 2006-2011) are closer to each other, denoting that most of the recent articles have received relatively same number of citations, due to short time span as it takes time for a (good) article to have an impact in the area. The three earliest articles (two in 1991 and one in 1992) have received relatively low citations. The article with the highest number of citations is a 2001 IEEE TSE article by Memon *et al.* titled *'Hierarchical GUI Test Case Generation Using Automated Planning'* [70] and has received 204 citations.

**RQ 2.5:** *What were the venues with highest articles count?* Figure 13 shows a count of articles from the top twenty venues, which contributed 80 articles. The annual International Workshop on TESTing Techniques & Experimentation Benchmarks for Event-Driven Software (TESTBEDS) is a relatively new venue, started in 2009. Since the venue has the specific focus on testing GUI and event-driven software, it has published
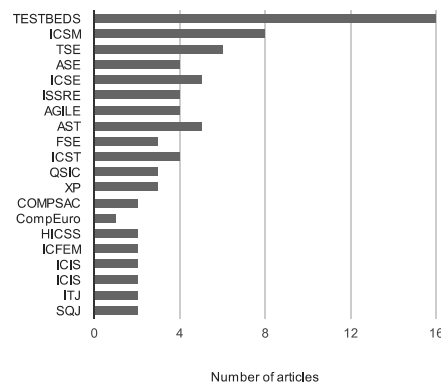


Figure 13: Data for **RQ 2.5:** Top Twenty Venues

the largest number, 16, of articles during 2009-2011. The International Conference on Software Maintenance (ICSM) with 8 and IEEE Transactions on Software Engineering (TSE) with 6 articles follow.

**RQ 2.6:** *What were the venues with highest citation count?* Figure 14 shows that the top three cited venues are (1) IEEE TSE, (2) ACM SIGSOFT Symposium on the Foundations of Software (FSE) (3) International Symposium on Software Reliability Engineering (ISSRE). Some venues such as FSE did not publish many GUI testing articles (3). However, those articles have received a large number of citations (349). The correlation between the number of articles in each venue versus the total number of citations to those articles was 0.46 (thus, not
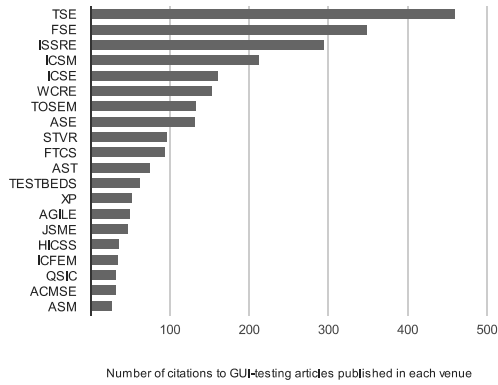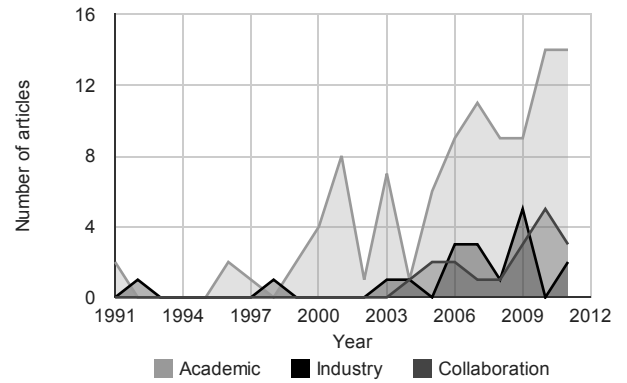
Figure 14: Data for **RQ 2.6**: Venues Most Cited



Figure 16: Data for **RQ 2.8**: Author Affiliation Trend

strong).

**RQ 2.7:** *Who are the authors with maximum articles?* As Figure 15 shows, Atif Memon (University of Maryland) stands first
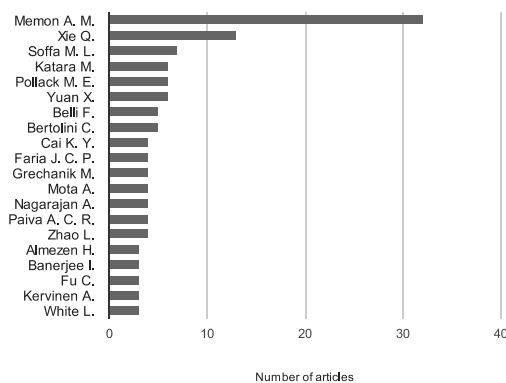


Figure 15: Data for **RQ 2.7**: Top 20 Authors

with 32 articles. The second and third highest ranking authors are Qing Xie (Accenture Tech Labs) and Mary Lou Soffa (University of Virginia) with 13 and 7 articles, respectively.

**RQ 2.8:** *What are the author affiliations, i.e., do they belong to academia or industry?* We classify the articles as coming from one of the following three categories based on the authors' affiliations: *academia*, *industry*, and *collaboration* (for articles whose authors come from both academics and industry). 73.52%, 13.23%, and 13.23% of the articles have been published by academics only, by industrial practitioners only, and with a collaboration between academic and industrial practitioners, respectively. The trend in each category over the years were tracked to see how many articles were written by academics or practitioners in different years. The results are shown in Figure 16. There is a steady rise in the number of articles published from academia and industry in recent years. Also the number of collaborative articles between academics and practitioners has been on the rise.

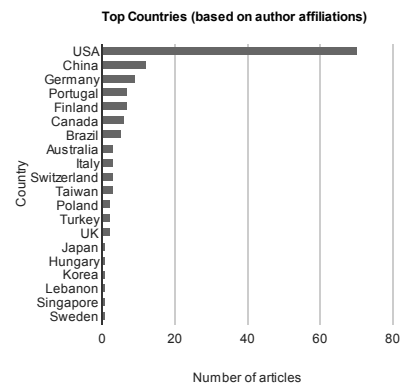**RQ 2.9:** *Which countries have produced the most articles?* To



Figure 17: Data for **RQ 2.9**: Top Author Countries

rank countries based on number of articles published, the country of the residence of the authors was extracted. If a article had several authors from several countries, one credit for each country was assigned.

The results are shown in Figure 17. The American researches have authored or co-authored 51.47% (70 of the 136) articles in the pool. Authors from China and Germany (with 12 and 9 articles, respectively) stand in the second and third ranks. Only 20 countries of the world have contributed to the GUI testing body of knowledge. International collaboration among the GUI testing researchers is quite under-developed as only 7 of the 136 articles were collaborations across two or more countries. Most of the remaining articles were written by researchers from one country.

## 8. Map Limitations & Future Directions

It is typical for research articles to state the limitations of the work and guidance for continuing research in the area. The research questions **RQ 3.*** are addressed by classifying the *reported* limitations and future directions from the articles.

**RQ 3.1:** *What limitations have been reported?* Many of the articles explicitly stated limitations of the work. The limitations were broadly categorized as follows:
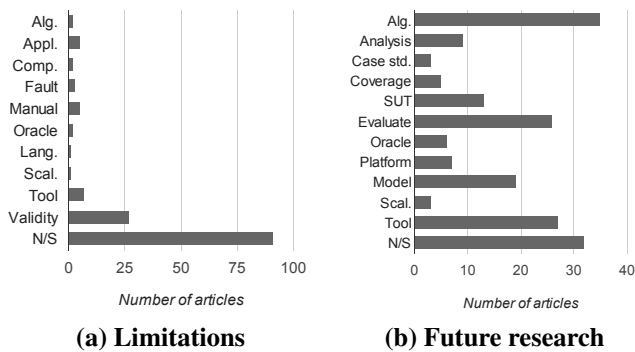
11

**(a) Limitations**  **(b) Future research**

Figure 18: Data for **RQ 3.***

---

• *algorithm:* The techniques or algorithms presented has known limitations - for example, an algorithm might not handle loops well [37].

• *applicability:* Limitations on usability under different environments - for example, a tool or algorithm may be specific to AWT based applications [103].

• *manual:* Manual steps are used in the experiments which may be limit the usability of the method. Manual steps may also affect the quality of the experiment or technique - for example, manual effort may be required to maintain a model [74].

• *oracle:* The oracle used for experiments may be limited in its capabilities at detecting all faults - for example, an oracle might be limited to detecting SUT crashes or exceptions [15], as opposed to comparing GUI states.

• *fault:* Limitations on the ability to detect all kinds of faults or may detect false defects - for example, a tool might not handle unexpected conditions well [22].

• *scalability:* Does not scale well to large GUIs - for example, time taken to execute the algorithm may increase super-linearly with large GUIs [41].

• *tool:* There is some known tool limitation or obvious missing features in the tools used or proposed - for example, a tool may handle only certain types of GUI events [8].

• *validity:* Experimental results are subject to internal or external validity [7, 17].

Out of the 136 articles, 45 articles reported one or more limitation of the research work. The extracted information is shown in Figure 18(a). This figure helps us understand the kind of limitations of the research work that were noted by the authors. The x-axis shows the number of articles in each category, the y-axis enumerates each category. The most common limitation is *validity*.

**RQ 3.2:** *What lessons learned are reported?* Only a small number of authors explicitly reported the lessons learned from their studies. Lesson learned were reported in only 11.76% (16/136) of all the articles. Lessons learned varied from author to author. They largely depend on individual research and study context. Hence, we conducted a qualitative analysis, instead of a quantitative analysis. It is important to note that they should be interpreted within the context of the studies.

Depending on the proposed testing techniques, the research lessons particularly associated with these techniques were reported. For example, in some cases, the authors who focus on model based testing where the model is created by hand, noted that in their approaches, a large amount effort would be spent on model creation [78]. In some other cases, the authors who used automated reverse engineered model based techniques, concluded that most of the tester's effort would be spent on test maintenance since the model is automatically created [98, 60]. The model in those techniques can be obtained at a low cost.

Similarly, the experimentation environment has also influenced the authors' suggestions. Some authors with limited computation resources suggested that more research effort should be spent on test selection [100], test prioritization [95] and test refactoring [30] to reduce the number of test cases to execute. However, some other authors with rich computation resources suggested that future research should focus on large scale studies [86].

**RQ 3.3:** *What are the trends in the area?* A widespread use of Java based SUTs and tools appears common. A notable development is the emergence of GUI testing work on mobile platforms during this period–8 article [8, 18, 15, 16, 17, 50, 51, 57]), compared to only 1 article in the period 1991-2007 [54].

Another notable trend is a shift from small unit script testing to large scale automated system testing. Several large scale empirical studies have been enabled [7, 10, 104]. thanks to the availability of automation tools and inexpensive computation resources.

**RQ 3.4:** *What future research directions are being suggested?* GUI testing is a relatively new research area in software engineering. Most of the articles provided guidance for continuing research, which may be broadly classified into the following categories:

• *algorithmic:* Extend existing algorithms or develop new ones - for example, extend the algorithm to handle potentially large number of execution paths [23].

• *analysis:* Further analyze results or techniques, further investigation based on results from the given study - for example, investigate interaction of different GUI components with CIS [95].

• *coverage:* Coverage techniques presented in the article can be further improve or evaluated. The coverage technique may be applicable for either code, GUI or model coverage - for example, develop new coverage criteria [108].

• *evaluate:* Evaluate the proposed methods, and techniques further, extend investigation based on existing results - for example, conduct more controlled experiments [15].

• *platform:* Extend the implementation for other platforms, *e.g.,* web and mobile [41].

• *model:* Improve or analyze the model presented in the article - for example, automatic generation of a model [76].

• *scalability:* Scale the proposed algorithms to larger systems, reduce computation cost - for example, scaling the algorithm to handle larger GUIs while improving execution performance [38].

• *SUT:* Evaluate the proposed techniques with a more SUTs - for example, use complex SUTs for evaluation [77].

12

• *tool:* Extend or add new capability or features to tools discusses in the article - for example, improve a tool to support better pattern matching and have better recovery from errors [27].

The future directions of research stated in the articles were extracted. Figure 18(b) shows this data. This figure helps us understand what guidance has been provided by researchers. Although this data contains *future* directions dating back to the year 1991, it helps us understand the thoughts of researchers during this period and what they perceived as missing pieces at the time their work was published.

In Figure 18(b) the x-axis shows the number of articles in each category, the y-axis enumerates each category. It can be seen that improving algorithms (35 articles) as been perceived as the area requiring them most work. Improving and developing better GUI testing tools has also been perceived as an area requiring work (27 articles).

## 9. Conclusions

This SM is the most comprehensive mapping of articles in the area of *GUI Testing*. A total of 230 articles, from the years 1991–2011, were collected and studied, from which 136 articles were included in the SM. Our findings indicate that most researchers work on developing new testing techniques or improving existing ones. Few articles express opinion about the state of the art in GUI testing. There is a large focus on model-based testing with models such as FSM, EFG and UML. There has been increased collaboration between academia and industry. However, no study has yet compared the state-of-the-art in GUI testing between academic and industrial tools and techniques.

An important result of this SM is that not all articles include information that is sought for secondary studies. We recommend that researchers working on GUI testing consider providing information in their articles using our maps as guides. In the future, we will continue to maintain an online repository [94] of GUI testing articles. We intend to continue analyzing the repository to create a systematic literature review (SLR).

## 10. Acknowledgements

## References

[1] GUITAR - A GUI Testing frAmewoRk. http://guitar.sourceforge.net.

[2] N. Abdallah and S. Ramakrishnan. Automated Stress Testing of Windows Mobile GUI Applications. In *Internation Symposium on Software Reliability Engineering*, 2009.

[3] W. Afzal, R. Torkar, and R. Feldt. A systematic mapping study on non-functional search-based software testing. In *20th International Conference on Software Engineering and Knowledge Engineering (SEKE 2008)*, 2008.

[4] W. Afzal, R. Torkar, and R. Feldt. A systematic review of search-based testing for non-functional system properties. *Inf. Softw. Technol.*, 51:957–976, June 2009.

[5] S. Ali, L. C. Briand, H. Hemmati, and R. K. Panesar-Walawege. A systematic review of the application and empirical investigation of search-based test case generation. *IEEE Trans. Softw. Eng.*, 36:742–762, November 2010.

[6] M. Alles, D. Crosby, C. Erickson, B. Harleton, M. Marsiglia, G. Pattison, and C. Stienstra. Presenter First: Organizing Complex GUI Applications for Test-Driven Development. In *Proceedings of the conference on AGILE 2006*, pages 276–288, 2006.

[7] D. Amalfitano, A. R. Fasolino, and P. Tramontana. Rich Internet Application Testing Using Execution Trace Data. In *Conference on Software Testing, Verification, and Validation Workshops*, pages 274–283, 2010.

[8] D. Amalfitano, A. R. Fasolino, and P. Tramontana. A GUI Crawling-Based Technique for Android Mobile Application Testing. In *Proceedings of the 2011 IEEE Fourth International Conference on Software Testing, Verification and Validation Workshops*, ICSTW '11, pages 252–261. IEEE Computer Society, 2011.

[9] O. E. Ariss, D. Xu, S. Dandey, B. Vender, P. McClean, and B. Slator. A Systematic Capture and Replay Strategy for Testing Complex GUI Based Java Applications. In *Conference on Information Technology*, pages 1038–1043, 2010.

[10] S. Arlt, C. Bertolini, and M. Schäf. Behind the Scenes: An Approach to Incorporate Context in GUI Test Case Generation. In *Proceedings of the 2011 IEEE Fourth International Conference on Software Testing, Verification and Validation Workshops*, ICSTW '11, pages 222–231. IEEE Computer Society, 2011.

[11] L. Baresi and M. Young. Test Oracles. Technical Report CIS-TR-01-02, University of Oregon, Dept. of Computer and Information Science, Eugene, Oregon, U.S.A., August 2001.

[12] Z. A. Barmi, A. H. Ebrahimi, and R. Feldt. Alignment of requirements specification and testing: A systematic mapping study. In *Proceedings of the 2011 IEEE Fourth International Conference on Software Testing, Verification and Validation Workshops*, ICSTW '11, pages 476–485, 2011.

[13] V. Basili, G. Caldiera, and H. Rombach. *Encyclopedia of Software Engineering*, chapter Goal Question Metric Approach, pages 528–532. John Wiley & Sons, Inc., 1994.

[14] F. Belli. Finite-State Testing and Analysis of Graphical User Interfaces. In *Symposium on Software Reliability Engineering*, page 34, 2001.

[15] C. Bertolini and A. Mota. Using Probabilistic Model Checking to Evaluate GUI Testing Techniques. In *Conference on Software Engineering and Formal Methods*, pages 115–124, 2009.

[16] C. Bertolini and A. Mota. A Framework for GUI Testing based on Use Case Design. In *Conference on Software Testing, Verification, and Validation Workshops*, pages 252–259, 2010.

[17] C. Bertolini, A. Mota, E. Aranha, and C. Ferraz. GUI Testing Techniques Evaluation by Designed Experiments. In *Conference on Software Testing, Verification and Validation*, pages 235–244, 2010.

[18] C. Bertolini, G. Peres, M. Amorim, and A. Mota. An Empirical Evaluation of Automated Black-Box Testing Techniques for Crashing GUIs. In *Software Testing Verification and Validation*, pages 21–30, 2009.

[19] R. V. Binder. Testing object-oriented software: a survey. In *Proceedings of the Tools-23: Technology of Object-Oriented Languages and Systems*, pages 374–, 1997.

[20] D. Budgen, M. Turner, P. Brereton, and B. Kitchenham. Using Mapping Studies in Software Engineering. In *Proceedings of PPIG 2008*, pages 195–204. Lancaster University, 2008.

[21] K.-Y. Cai, L. Zhao, H. Hu, and C.-H. Jiang. On the Test Case Definition for GUI Testing. In *Conference on Quality Software*, pages 19–28, 2005.

[22] T.-H. Chang, T. Yeh, and R. C. Miller. GUI Testing Using Computer Vision. In *Conference on Human factors in computing systems*, pages 1535–1544, 2010.

[23] J. Chen and S. Subramaniam. Specification-based Testing for GUI-based Applications. *Software Quality Journal*, 10(2):205–224, 2002.

[24] W.-K. Chen, T.-H. Tsai, and H.-H. Chao. Integration of Specification-Based and CR-Based Approaches for GUI Testing. In *Conference on Advanced Information Networking and Applications*, pages 967–972, 2005.

[25] V. Chinnapongsea, I. Lee, O. Sokolsky, S. Wang, and P. L. Jones. Model-Based Testing of GUI-Driven Applications. In *Workshop on Software Technologies for Embedded and Ubiquitous Systems*, pages 203–214, 2009.

[26] K. Conroy, M. Grechanik, M. Hellige, E. Liongosari, and Q. Xie. Automatic Test Generation from GUI Applications for Testing Web Services. In *Conference on Software Maintenance*, pages 345–354, 2007.

[27] M. Cunha, A. Paiva, H. Ferreira, and R. Abreu. PETTool: A pattern-based GUI testing tool. In *International Conference on Software Technology and Engineering*, pages 202–206, 2010.

[28] P. A. da Mota Silveira Neto, I. d. Carmo Machado, J. D. McGregor, E. S. de Almeida, and S. R. de Lemos Meira. A systematic mapping study of software product lines testing. *Inf. Softw. Technol.*, 53:407–423, 2011.

[29] T. Daboczi, I. Kollar, G. Simon, and T. Megyeri. Automatic Testing of Graphical User Interfaces. In *Instrumentation and Measurement Technology Conference*, pages 441–445, 2003.

[30] B. Daniel, Q. Luo, M. Mirzaaghaei, D. Dig, D. Marinov, and M. Pezzè. Automated GUI refactoring and test script repair. In *Proceedings of the First International Workshop on End-to-End Test Script Engineering*, ETSE '11, pages 38–41, New York, NY, USA, 2011. ACM.

[31] A. Derezinska and T. Malek. Unified Automatic Testing of a GUI Applications' Family on an Example of RSS Aggregators. In *Multiconference on Computer Science and Information Technology*, pages 549–559, 2006.

[32] A. Derezinska and T. Malek. Experiences in Testing Automation of a Family of Functional-and GUI-similar Programs. *Journal of Computer Science & Applications*, 4(1):13–26, 2007.

[33] M. B. Dwyer, V. Carr, and L. Hines. Model Checking Graphical User Interfaces Using Abstractions. In *ESEC / SIGSOFT FSE*, pages 244–261, 1997.

[34] L. Feng and S. Zhuang. Action-driven Automation Test Framework for Graphical User Interface (GUI) Software Testing. In *Autotestcon*, pages 22 – 27, 2007.

[35] A. Fernandez, E. Insfran, and S. Abrah£o. Usability evaluation methods for the web: A systematic mapping study. *Information and Software Technology*, 53(8):789 – 817, 2011.

[36] S. Ganov, C. Killmar, S. Khurshid, and D. E. Perry. Event Listener Analysis and Symbolic Execution for Testing GUI Applications. *Formal methods and software engineering*, 5885(1):69–87, 2009.

[37] S. Ganov, C. Kilmar, S. Khurshid, and D. Perry. Test Generation for Graphical User Interfaces Based on Symbolic Execution. In *Proceedings of the International Workshop on Automation of Software Test*, 2008.

[38] R. Gove and J. Faytong. Identifying Infeasible GUI Test Cases Using Support Vector Machines and Induced Grammars. In *Proceedings of the 2011 IEEE Fourth International Conference on Software Testing, Verification and Validation Workshops*, pages 202–211, 2011.

[39] M. Grechanik, D. S. Batory, and D. E. Perry. Integrating and Reusing GUI-Driven Applications. In *ICSR*, pages 1–16, 2002.

[40] M. Grechanik, Q. Xie, and C. Fu. Experimental Assessment of Manual Versus Tool-based Maintenance of GUI-Directed Test Scripts. *Conference on Software Maintenance*, pages 9–18, 2009.

[41] M. Grechanik, Q. Xie, and C. Fu. Maintaining and Evolving GUI-directed Test Scripts. In *Conference on Software Engineering*, pages 408–418, 2009.

[42] M. Grindal, J. Offutt, and S. F. Andler. Combination testing strategies: A survey. *Software Testing, Verification, and Reliability*, 15:167–199, 2005.

[43] M. J. Harrold. Testing: a roadmap. In *Proceedings of the Conference on The Future of Software Engineering*, ICSE '00, pages 61–72, New York, NY, USA, 2000. ACM.

[44] S. Herbold, J. Grabowski, S. Waack, and U. Bünting. Improved bug reporting and reproduction through non-intrusive gui usage monitoring and automated replaying. In *Proceedings of the 2011 IEEE Fourth International Conference on Software Testing, Verification and Validation Workshops*, ICSTW '11, pages 232–241, 2011.

[45] A. Holmes and M. Kellogg. Automating Functional Tests Using Selenium. In *agile Conference*, pages 270–275, 2006.

[46] Y. Hou, R. Chen, and Z. Du. Automated GUI Testing for J2ME Software Based on FSM. In *Conference on Scalable Computing and Communications*, pages 341–346, 2009.

[47] C. Hu and I. Neamtiu. Automating gui testing for android applications. In *Proceedings of the 6th International Workshop on Automation of Software Test*, pages 77–83, 2011.

[48] C. Hu and I. Neamtiu. A gui bug finding framework for android applications. In *Proceedings of the 2011 ACM Symposium on Applied Computing*, SAC '11, pages 1490–1491. ACM, 2011.

[49] Z. Hui, R. Chen, S. Huang, and B. Hu. Gui regression testing based on function-diagram. In *Intelligent Computing and Intelligent Systems (ICIS), 2010 IEEE International Conference on*, volume 2, pages 559 –563, oct. 2010.

[50] A. Jaaskelainen, M. Katara, A. Kervinen, M. Maunumaa, T. Paakkonen, T. Takala, and H. Virtanen. Automatic GUI test generation for smartphone applications - an evaluation. In *Conference on Software Engineering*, pages 112–122, 2009.

[51] A. Jaaskelainen, A. Kervinen, and M. Katara. Creating a Test Model Library for GUI Testing of Smartphone Applications. In *Conference on Quality Software*, pages 276–282, 2008.

[52] Y. Jia and M. Harman. An analysis and survey of the development of mutation testing. *IEEE Transactions on Software Engineering*, 2008:1–32, 2010.

[53] N. Juristo, A. M. Moreno, and S. Vegas. Reviewing 25 years of testing technique experiments. *Empirical Softw. Engg.*, 9:7–44, 2004.

[54] A. Kervinen, M. Maunumaa, T. Paakkonen, and M. Katara. Model-Based Testing Through a GUI. *Formal approaches to software testing*, 3997(1):16–31, 2006.

[55] B. Kitchenham and S. Charters. Guidelines for performing systematic literature reviews in software engineering. *Version*, 2(EBSE 2007-001):200701, 2007.

[56] B. A. Kitchenham, D. Budgen, and O. P. Brereton. Using mapping studies as the basis for further research ¢ a participant-observer case study. *Information and Software Technology*, 53(6):638 – 651, 2011.

[57] O.-H. Kwon and S.-M. Hwang. Mobile GUI Testing Tool based on Image Flow. In *Conference on Computer and Information Science*, pages 508–512, 2008.

[58] P. Li, T. Huynh, M. Reformat, and J. Miller. A Practical Approach to Testing GUI Systems. *Empirical software engineering*, 12(4):331–357, 2007.

[59] R. Lo, R. Webby, and R. Jeffery. Sizing and Estimating the Coding and Unit Testing Effort for GUI Systems. In *Software Metrics Symposium*, page 166, 1996.

[60] C. Lowell and J. Stell-Smith. Successful Automation of GUI Driven Acceptance Testing. *Extreme programming and Agile processes in software engineering*, 2675(1):1011–1012, 2003.

[61] K. Magel and I. Alsmadi. GUI Structural Metrics and Testability Testing. In *Conference on Software Engineering and Applications*, pages 91–95, 2007.

[62] S. A. M. Mark Harman and Y. Zhang. Search based software engineering: A comprehensive analysis and review of trends techniques and applications. Technical Report TR-09-03, Department of Computer Science, King's College London, April 2009.

[63] S. McConnell. Daily Build and Smoke Test. *IEEE Software*, 13(4):143–144, 1996.

[64] C. McMahon. History of a Large Test Automation Project Using Selenium. In *Proceedings of the 2009 Agile Conference*, pages 363–368, 2009.

[65] P. McMinn. Search-based software test data generation: a survey: Research articles. *Softw. Test. Verif. Reliab.*, 14:105–156, June 2004.

[66] A. M. Memon. Automatically Repairing Event Sequence-based GUI Test Suites for Regression Testing. *ACM Transactions on Software Engineering and Methodology*, 18(2):1–36, 2008.

[67] A. M. Memon and B. N. Nguyen. Advances in automated model-based system testing of software applications with a GUI front-end. In M. V. Zelkowitz, editor, *Advances in Computers*, volume 80, pages nnn–nnn. Academic Press, 2010.

[68] A. M. Memon, M. E. Pollack, and M. L. Soffa. Automated Test Oracles for GUIs. *ACM SIGSOFT Software Engineering Notes*, 25(6):30–39,

2000.

[69] A. M. Memon, M. E. Pollack, and M. L. Soffa. Plan Generation for GUI Testing. In *Conference on Artificial Intelligence Planning and Scheduling*, pages 226–235, 2000.

[70] A. M. Memon, M. E. Pollack, and M. L. Soffa. Hierarchical GUI Test Case Generation Using Automated Planning. *IEEE Transactions on Software Engineering*, 27(2):144–155, 2001.

[71] A. M. Memon, M. L. Soffa, and M. E. Pollack. Coverage Criteria for GUI Testing. In *Software Engineering conference held jointly with ACM SIGSOFT symposium on Foundations of software engineering*, pages 256–267, 2001.

[72] B. A. Myers. User interface software tools. *ACM Trans. Comput.-Hum. Interact.*, 2:64–103, March 1995.

[73] M. Navarro, P. Luis, S. Ruiz, D. M. Perez, and Gregorio. A Proposal for Automatic Testing of GUIs Based on Annotated Use Cases. *Advances in Software Engineering*, 2010(1):1–8, 2010.

[74] D. H. Nguyen, P. Strooper, and J. G. Suess. Model-Based Testing of Multiple GUI Variants Using the GUI Test Generator. In *Workshop on Automation of Software Test*, pages 24–30, 2010.

[75] H. Okada and T. Asahi. GUITESTER: A Log-Based Usability Testing Tool for Graphical User Interfaces. *IEICE Transactions on Information and Systems*, 82:1030–1041, 1999.

[76] A. C. Paiva, J. C. Faria, N. Tillmann, and R. A. Vidal. A Model-to-Implementation Mapping Tool for Automated Model-Based GUI Testing. *Formal methods and software engineering*, 3785(1):450–464, 2005.

[77] A. C. R. Paiva, J. C. P. Faria, and P. M. C. Mendes. Reverse Engineered Formal Models for GUI Testing. *Formal methods for industrial critical systems*, 4916(1):218–233, 2008.

[78] A. C. R. Paiva, N. Tillmann, J. C. P. Faria, and R. F. A. M. Vidal. Modeling and Testing Hierarchical GUIs. In *Workshop on Abstract State Machines*, 2005.

[79] M. Palacios, J. García-Fanjul, and J. Tuya. Testing in service oriented architectures with dynamic binding: A mapping study. *Inf. Softw. Technol.*, 53:171–189, 2011.

[80] R. M. Patton and G. H. Walton. An Automated Testing Perspective of Graphical User Interfaces. In *The Interservice/Industry Training, Simulation & Education Conference*, 2003.

[81] K. Petersen, R. Feldt, S. Mujtaba, and M. Mattsson. Systematic Mapping Studies in Software Engineering. *12th International Conference on Evaluation and Assessment in Software Engineering*, 17(1):1–10, 2007.

[82] J. Portillo-Rodriguez, A. Vizcaino, M. Piattini, and S. Beecham. Tools used in global software engineering: a systematic mapping review. *Information and Software Technology*, 2012.

[83] C. S. Pǎsǎreanu and W. Visser. A survey of new trends in symbolic execution for software testing and analysis. *Int. J. Softw. Tools Technol. Transf.*, 11(4):339–353, Oct. 2009.

[84] B. Robinson and P. Brooks. An Initial Study of Customer-Reported GUI Defects. In *Conference on Software Testing, Verification, and Validation Workshops*, pages 267–274, 2009.

[85] M. Safoutin, C. Atman, R. Adams, T. Rutar, J. Kramlich, and J. Fridley. A design attribute framework for course planning and learning assessment. *Education, IEEE Transactions on*, 43(2):188 –199, may 2000.

[86] Y. Shewchuk and V. Garousi. Experience with Maintenance of a Functional GUI Test Suite using IBM Rational Functional Tester. In *Proceedings of the International Conference on Software Engineering and Knowledge Engineering*, pages 489–494, 2010.

[87] B. Shneiderman and C. Plaisant. *Designing the User Interface - Strategies for Effective Human-Computer Interaction (5. ed.)*. Addison-Wesley, 2010.

[88] J. Strecker and A. Memon. Relationships Between Test Suites, Faults, and Fault Detection in GUI Testing. In *Conference on Software Testing, Verification, and Validation*, pages 12–21, 2008.

[89] J. Takahashi. An Automated Oracle for Verifying GUI Objects. *ACM SIGSOFT Software Engineering Notes*, 26(4):83–88, 2001.

[90] F. M. Theodore D. Hellmann, Ali Hosseini-Khayat. *Agile Interaction Design and Test-Driven Development of User Interfaces - A Literature Review*. Number 9. Springer, 2010.

[91] Y. Tsujino. A verification method for some GUI dialogue properties. *Systems and Computers in Japan*, pages 38–46, 2000.

[92] `http://crestweb.cs.ucl.ac.uk/resources/sbse_repository/`.

[93] `http://www.cs.umd.edu/~atif/testbeds/testbeds2011.htm`.

[94] `http://www.softqual.ucalgary.ca/projects/2012/GUI_SM/`.

[95] L. White and H. Almezen. Generating Test Cases for GUI Responsibilities Using Complete Interaction Sequences. In *Symposium on Software Reliability Engineering*, page 110, 2000.

[96] L. J. White. Regression Testing of GUI Event Interactions. In *Conference on Software Maintenance*, pages 350–358, 1996.

[97] Q. Xie and A. Memon. Rapid "Crash Testing" for Continuously Evolving GUI-Based Software Applications. In *Conference on Software Maintenance*, pages 473–482, 2005.

[98] Q. Xie and A. M. Memon. Designing and Comparing Automated Test Oracles for GUI-based Software Applications. *ACM Transactions on Software Engineering and Methodology*, 16(1):1–36, 2007.

[99] Q. Xie and A. M. Memon. Using a Pilot Study to Derive a GUI Model for Automated Testing. *ACM Transactions on Software Engineering and Methodology*, 18(2):1–33, 2008.

[100] M. Ye, B. Feng, Y. Lin, and L. Zhu. Neural Networks Based Test Cases Selection Strategy for GUI Testing. In *Congress on Intelligent Control and Automation*, pages 5773–5776, 2006.

[101] S. Yip and D. Robson. Applying Formal Specification and Functional Testing to Graphical User Interfaces. In *Advanced Computer Technology, Reliable Systems and Applications European Computer Conference*, pages 557 – 561, 1991.

[102] S. Yip and D. Robson. Graphical User Interfaces Validation: a Problem Analysis and a Strategy to Solution. In *Conference on System Sciences*, 1991.

[103] X. Yuan, M. B. Cohen, and A. M. Memon. GUI Interaction Testing: Incorporating Event Context. In *IEEE Transactions on Software Engineering*, 2010.

[104] X. Yuan and A. M. Memon. Using GUI Run-Time State as Feedback to Generate Test Cases. In *Proceedings of the 29th international conference on Software Engineering*, ICSE '07, pages 396–405, Washington, DC, USA, 2007. IEEE Computer Society.

[105] X. Yuan and A. M. Memon. Generating Event Sequence-Based Test Cases Using GUI Runtime State Feedback. *IEEE Trans. Softw. Eng.*, 36:81–95, January 2010.

[106] X. Yuan and A. M. Memon. Iterative Execution-feedback Model-directed GUI Testing. *Information and Software Technology*, 52(5):559–575, 2010.

[107] H. Zhang, M. A. Babar, and P. Tell. Identifying relevant studies in software engineering. *Inf. Softw. Technol.*, 53:625–637, 2011.

[108] L. Zhao and K.-Y. Cai. Event Handler-Based Coverage for GUI Testing. In *Conference on Quality Software*, pages 326–331, 2010.