# Flycatcher

Automatic unit test generation for JavaScript

*Interim Report*

**Author:** Jerome DE LAFARGUE
**Supervisor:** Susan EISENBACH
**Technical Supervisor:** Tristan ALLWOOD

DEPARTMENT OF COMPUTING
IMPERIAL COLLEGE LONDON

# Contents

# Chapter 1

# Introduction

Software testing is a cornerstone of software engineering — one of the most common and effective ways to verify software quality and an effort that accounts for at least 50% of software development time [32]. With the fast-paced growth of the software industry, comes the need to test larger and more complex software on an unprecedented scale. Moreover, as software becomes increasingly ubiquitous, it is held to the highest standards of reliability and correctness, which further justifies testing it in a rigorous and exhaustive manner.

As a result, many attempts have been made to automate the testing effort, so that programs can be systematically and seamlessly tested, without requiring laborious, costly and error-prone manual input. The consequences of automated testing are very appealing: it reduces software maintenance and development costs, while increasing the robustness and ultimate quality of the software. Despite the fact that this area of research has taken time to develop, due to the intrinsic complexities of automatic test generation, it has now seemingly reached a stage where it can start to make a meaningful impact on software testing practice.

Decades of research have been devoted to automatic test generation for static languages and a multitude of tools have been developed. As the research area matures, it is arriving to a point where its techniques are no longer simply applicable to restricted programming language subsets or limited programs. Indeed, companies such as Microsoft employ automatic test generation tools on a regular basis to verify their software [26]. Yet, until very recently, dynamic programming languages had been left out of the equation — but their increasing popularity and a renewed interest in them prompts the need to start including them in the automatic testing research effort.

One such programming language that has been growing in popularity in the past few years is JavaScript, with new frameworks and libraries being released for it frequently. Software libraries that have gained wide acceptance, like NODE.JS[1] which supports the writing of highly-scalable internet applications, seem to confirm JavaScript's transition from a purely client-side browser language to a an all-

---

[1]http://nodejs.org/

purpose one — at least for some. In recent studies [17], JavaScript appears amongst the most used programming languages in the world today. In other words, it seems that JavaScript is here to stay, at least for some time, and it makes sense to devote time to it.

Various test data generation methods exist from the simple and imprecise random test generation to hugely complex and elaborate systems that combine static and dynamic analysis to provide strong software verification. Since much of the literature on automatic test generation focuses on static languages and numerical data types, many of the techniques found are not feasible or applicable to automatic test generation for JavaScript. However, some are more appropriate or feasible than others for our objective — to generate unit test suites for JavaScript objects — and it is those methods that we will explore during the course of this project.

This project is challenging for several reasons. First of all, it requires automatic test data generation, which is an undecidable task in itself []

... Thankfully, ... Secondly ...

The major contribution that this project hopes to make to the field is to extend the limited amount of work done regarding dynamic languages by proposing a tool for automatic unit test generation for JavaScript programs. As well as offering a potential tool for JavaScript developers, we hope that our work will be able to offer new insights into automatic test generation for a dynamic language and benefit future research in that direction.

# Chapter 2

# Background

In this chapter we will start by giving an overview of software testing, with particular emphasis on aspects of it that are relevant to this project. We will then take a look at the state of the art in automatic test data generation (ATDG) in order to understand the approach that will be used for Flycatcher. Because the tests will be object-oriented, method call sequences also need to be generated for the test cases and we will look at the state of the art for doing that too. Finally we will explain why we chose JavaScript as our target language and describe features of it that are of interest to us for this project.

## 2.1 Dynamic software testing

### 2.1.1 Overview

We can define the activity of dynamic testing as testing that requires execution of the software with test data as input [18] and characterise it with respect to three parameters namely the amount of knowledge assumed by the tester, the target of the tests and the stage of development at which they are executed. The amount of knowledge of the software under test can be divided into three categories, structural (white-box testing) testing, functional (black-box testing) and a hybrid of the two (grey-box testing). The target of the tests refers to their granularity, from testing specific units of code (unit testing) to an entire integrated system (system testing). The stage at which the tests are undertaken determines whether they are regression tests, alpha-tests, beta-tests, acceptance tests *etc*. With Flycatcher we hope to generate suites of structural tests, focused at the unit level of object-oriented classes, most likely to perform incremental regression testing. Hence, structural testing, unit testing and regression testing will be described in more detail in this section.

### 2.1.2 Structural testing

The goal of structural testing is to test the internal workings [20] of an application in order to verify that it does not contain errors. While functional testing determines

whether the software provides the required functionality, structural testing tries to ensure that it does not crash under any circumstances, regardless of how it is called. It concerns *how* well the software operates, its structure, rather than *what* it can do, its function. As a result, the measure to determine good structural testing is the code covered during the testing process — code coverage. It gives us an idea of the amount of code that should be bug free. However, there are various types of code coverage criteria and the confidence that our code is bug free varies depending on which one is chosen.

**Test coverage**

Edvardsson lists the most cited criteria [7], from weakest to strongest:

- **Statement Coverage** Each statement must be executed at least once.

- **Branch[1]/Decision Coverage** Each branch condition must evaluate to true and false.

- **Condition/Predicate Coverage** Each clause within each branch condition must evaluate to true and false.

- **Multiple-condition Coverage** Each combination of truth values of each clause of each condition must be executed.

- **Path[1] Coverage** Each path in the control flow graph[1] must be traversed.

  The stronger criteria of condition, multiple-condition and path coverage are often infeasible to achieve for programs of more than moderate complexity, and thus branch coverage has been recognised as a basic measure for testing [7].

### 2.1.3   Unit testing

Unit testing consists in testing individual and independently testable units of source code [24]. Therefore, unit testing is made easier if the code is designed in a modular way. The nature of the units depends on the programming language and environment but is often a class or a function. As opposed to system tests which can be aimed at the client, unit tests are usually white-box tests. Although they do not guarantee that the overall software works as required, they give confidence in specific units of code and narrow down errors, helping the development process. In Flycatcher, the target unit will be a JavaScript prototype, which will be introduced later in this chapter.

---

[1]For an explanation of program analysis terminology such as branch, path and control flow graph see Appendix A.

### 2.1.4 Regression testing

Automatically generating structural unit tests can be of great use for regression testing. Regression testing aims to ensure that enhancing a piece of software does not introduce new faults [24]. The difficulty in testing this is that programmers do not always appreciate the extent of their changes. Hence, having a suite of unit tests with good structural coverage can reduce this problem by verifying the software in a systematic, unbiased way.

## 2.2 Automatic test data generation

### 2.2.1 Overview

As can be seen in Mahmood's systematic review of ATDG techniques [18], many classifications exist for ATDG techniques. For our purposes, the first distinction that we need to make is between white-box, black-box [28] and grey-box ATDG techniques, as for Flycatcher we are only interested in white-box testing. In the literature we found that white-box ATDG techniques are usually classified in two ways [18, 7, 32].

The first concerns the target selection stage of ATDG techniques: where either paths or individual nodes that contribute to the overall coverage criterion are successively selected from the control flow graph, so that test data that respectively traverses the path or reaches the node can be generated. When specific paths are targeted, the ATDG technique is known as *path-oriented* [7] whereas if a node is targeted then it is *goal-oriented*. When data is generated purely randomly *i.e.* there is no specific target, then as part of this classification the ATDG technique is simply *random*.

The other classification of white-box ATDG concerns the type of implementation: *static*, *dynamic* or a *hybrid* of the two [12, 20]. We will focus on the latter classification of structural testing as it governs our choice of implementation for Flycatcher. Moreover, the former concerns the path selection stage of ATDG and this step will be ignored in Flycatcher as it is in many recent ATDG techniques [32]. Figure 2.1 summarises what we believe is an intuitive characterisation of ATDG techniques with respect to this project and the one which will guide our choice of implementation. Many techniques can be found under each of the static, dynamic and hybrid implementation categories and we will only list the most noteworthy to us.

The choice of implementation for Flycatcher is dynamic random ATDG for our benchmark and dynamic search-based ATDG using genetic algorithms for our solution. The rest of this section will present in further detail the structural ATDG implementation categories we have chosen and the difficulties of ATDG for a dynamic language, so that we can understand the rationale behind this choice.
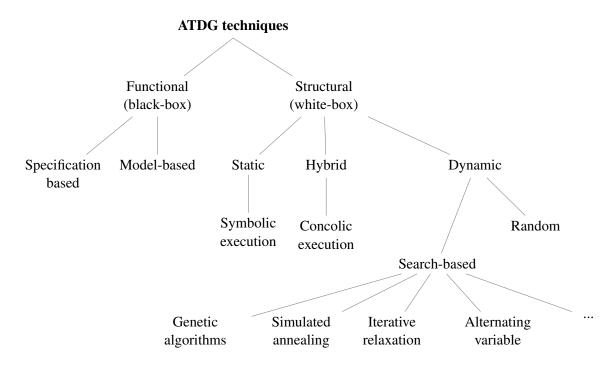
Figure 2.1: Overview of ATDG techniques

## 2.2.2 Static test data generation

Static structural test data generation is based on information available from the static analysis of the program, without requiring that the program is actually executed [20]. Static program analysis produces control flow information that can be used to select execution paths in order to try and achieve good coverage. The goal of ATDG is then to generate data that executes these paths.

Every time control flow branches, *e.g.* at `if` statements, there is a corresponding predicate or branch condition. These predicates can be collected along a path and conjoined to form the path predicate. By solving the path predicate in terms of the input variables, we can obtain test data that executes that path. However, in order to rewrite the path predicate in terms of the input variables we need to take into account the execution of the program. Hence, for generation of test data statically a technique called symbolic execution [14] is used.

Symbolic execution gathers constraints along a simulated execution of a program path, where symbolic variables are used instead of actual values, such that the final path predicate can be rewritten in terms of the input variables. Solving the resulting system of constraints then yields the data necessary for the traversal of the given path [13, 14]. There are a lot of technical difficulties associated with symbolic execution [7, 21, 20]:

- the presence of input variable dependent loops can lead to infinite execution

trees[2] as the loops can be executed any number of times

- array references become problematic if the indexes are not constants but variables, as is typically the case

- features such as pointers and dynamically-allocated objects that rely on execution are hard to analyse statically

- static analysis is not possible for function calls to precompiled modules or libraries

- if the path constraint is non-linear, solving it is an undecidable problem

- even if the path constraint is linear, solving it can lead to very high complexity

Although various static solutions have been proposed for these issues [29, 10, 25], they often dramatically increase the complexity of the ATDG process. As a result, tools purely based on symbolic execution can typically handle only subsets of programming languages and are not applicable in industry. A better trend that has developed in the past decade, is the combination of concrete and symbolic execution, which tackles most of the aforementioned issues [26] — we will cover this type of ATDG implementation in the subsection on hybrid ATDG. Due to the numerous problems posed by purely static ATDG, its weakness with dynamic types and constructs [7, 32] and the complexity of building a fully-fledged symbolic executor for a language [7, 12], we chose not to use static ATDG for the implementation of Flycatcher.

### 2.2.3   Hybrid test data generation

The hybrid approach to ATDG consists in combining symbolic and concrete execution, which is known as *concolic execution* [26]. In other words, hybrid analysis tools run programs on actual inputs, while collecting symbolic constraints in order to direct the search for new inputs. In doing so, they avoid the main weaknesses of the static approach, such as solving non-linear constraints or dealing with dynamic structures. This type of technique has been popular in recent years, mainly because it overcame the limitations that prevented static ATDG techniques from being applied to industry software. Notable tools that implement it are DART [9], CUTE [31], JPF-SE [2], PEX [33], EXE [5] and KLEE [4].

Yet, although it deals with some limitations of static ATDG, hybrid ATDG still requires static analysis of the source code under test, which is unfeasible for Flycatcher, given the dynamically typed, object-oriented nature of JavaScript.

---

[2]the execution paths followed during the symbolic execution of a procedure [14]

### 2.2.4 Dynamic test data generation

Dynamic structural test data generation is purely based on actual execution of the software. The program under test is run with, possibly randomly, selected input and feedback is collected at runtime regarding the chosen coverage objective [7]. The feedback is usually obtained through some form of instrumentation of the program that monitors the program flow. Inputs can be continually generated randomly, relying on probability to achieve the coverage objective — this is known as *random* test data generation and does not perform well in general [7]. On the other hand, inputs can be incrementally tuned based on the feedback (using different kinds of search methods) in order to satisfy the coverage objective — this is known as *search-based* test data generation [20], where the search-space is the control flow graph of a program. The main drawback of dynamic ATDG is that it is reliant on the speed of execution of a program and as the number of required executions to achieve satisfactory coverage may be high, this leads to an overall expensive process. Below we will present the random and search-based approaches in more detail, in order to understand which would be more suitable for Flycatcher.

#### Random approach

Random test data generation consists in producing inputs at random in the hope of achieving the chosen coverage criterion through probability. Although random test data generation is relatively simple to implement, it does not perform well in terms of coverage, as the chances of finding faults that are revealed by only a small percentage of program inputs are low [7]. In other words, it is difficult for it to exercise 'deep' features of a program that are exercised only through specific and unlikely paths. As a result, random ATDG only works well for straightforward programs. However, because it is the simplest ATDG technique and is considered to have the lowest acceptance rate [7], it is often used as a benchmark and is a suitable candidate for us to benchmark our Flycatcher application.

#### Search-based approach

Search-based test data generation uses heuristics to guide the generation of input data so that the inputs execute paths that contribute to the overall test coverage objective. This involves modelling the test coverage objective as a heuristic function or *objective function*, that evaluates the fitness of a chosen set of inputs with respect to a coverage objective. Based on those fitness values, many search techniques exist to find optimal inputs in order to achieve the desired coverage. Various objective functions exist and they are dependent on the ATDG method used. Some of the well-known search-based ATDG techniques are *alternating variable* (local search optimisation) [15, 8], *simulated annealing* [35, 34], *iterative relaxation* [11] and *genetic algorithms* [22, 23].

In 2008, Han and Kwon [12] conducted a robust empirical evaluation of search-based test data generation techniques, comparing iterative-relaxation, local search

optimisation and genetic algorithms. Genetic algorithms came out on top of the other techniques regarding the rate of coverage and generality, both essential characteristics when considering ATDG techniques. However, genetic algorithms proved expensive both in time and resources, but this can be improved upon [12]. Although simulated annealing is not part of this empirical evaluation, we can see from [22] that although it performs as well as genetic algorithms, it appears to be less efficient.

Additionally, genetic algorithms offer another advantage: as opposed to most other ATDG techniques, they do not require static analysis of the program under test [12]. This advantage is especially significant in our case, as to the best of our knowledge, tools that produce control and data flow analysis information for JavaScript source code are not available. However, this is not surprising since the dynamic, object-oriented nature of JavaScript makes this task very laborious.

Consequently, due to their effectiveness in terms of coverage and the advantages they possess with regard to dynamic languages, genetic algorithms appear to be the most suitable for our implementation of ATDG for Flycatcher. We will therefore go on to explain what genetic algorithms are and how they can be applied to test data generation.

### 2.2.5 Genetic algorithms

**Overview**

Genetic algorithms attempt to model the robustness and flexibility of natural selection in order to guide a search[23]. They start with a randomly initialised *population* of potential candidate solutions, called *individuals* or *chromosomes*. The population is iteratively recombined and mutated to evolve successive populations, known as *generations*. The recombination takes the 'fittest' parent solutions and 'breeds' them to produce new offspring. Their fitness is determined by a *fitness function* that evaluates how good a candidate solution is for a particular problem.

This process favours evolution towards fitter individuals, mimicking natural selection. The 'breeding' of two individuals usually involves a crossover operation which swaps their information at a randomly selected position. In the interest of maintaining diversification, a *mutation* phase usually occurs after the crossover, to introduce new genetic material into the search and avoid premature convergence on one area of the search space. The mutation operation randomly modifies some information of a selected individual, The process of iterative recombination and mutation, illustrated in Figure 2.2, is repeated until a specific termination criterion is fulfilled.

**Application to test data generation**

For the purposes of ATDG, the population is one of test cases. The test cases are evolved according to the genetic algorithm in order to satisfy the chosen coverage criteria [23]. In other words, the *fitness function* is based on the *objective function*,
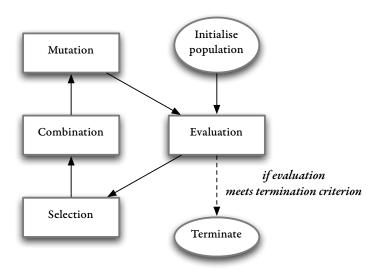
Figure 2.2: Flowchart of a genetic algorithm

the function that evaluates test data according to a certain coverage criterion — the more a candidate test case contributes to coverage, the fitter it is. Instead of using search heuristics to minimise our objective function, we use it to guide our natural selection process towards the test cases that will achieve the best coverage.

Regarding the objective function used in the context of genetic algorithms, different variants exist [20]. The *coverage-oriented approach* rewards individuals on the basis of *all* covered program structures. This type of objective function rewards coverage with respect to the overall goal *i.e.* the system is not required to pick a specific path or node and attempt to execute it. The other approaches are *structure-oriented* where, similarly to most other search-based ATDG methods, a path or node is selected from the data flow analysis and a separate search is undertaken for each selected structure. Once more, seeing that the former does not require static analysis of the program under test, it is the most attractive to us for the implementation of Flycatcher.

Finally, the population, individual representation, fitness evaluation and methods of crossover, mutation and selection are all highly dependent on the test data generation problem at hand, hence we will discuss them with regard to Flycatcher when we come to the implementation.

### 2.2.6 Challenges of dynamic object-oriented languages

Most of the research on ATDG concerns static programming languages [18] and it is only in the past few years that dynamic programming languages have sparked some interest in that field. A possible reason for this is that dynamic programming languages make ATDG harder by enabling features that allow programs to signif-

11

icantly change at runtime. These features can include modifying the type system, extending objects or adding new code during program execution. The challenges this type of behaviour introduces, and that Flycatcher will try and overcome, are listed below [6].

**Generating test data of the required type**

Given that method parameters do not have static types in dynamically typed languages, we do not know what arguments to pass to them. A potential solution to this is to use a method called *type inference* [27], which tries to infer the type of arguments from the way they are used inside the program. Although this method does not guarantee 100% precision, it is a good starting point for generating accurately typed test data. Mairhofer uses this technique for RUTEG [19], his search-based ATDG tool for Ruby, where the search for test data refines the initially inferred type, by discarding poor candidates.

**Generating object instances**

Sometimes the type of an object constructor or method parameters will be a complex type and this complicates the test data generation task even further. Generating well-formed object instances to use as arguments inside tests for a dynamically typed object-oriented language is problematic because there isn't a blueprint to construct them from. There is previous work on input data generation for dynamic data structures [15, 36, 30, 37], but all these approaches focus on statically typed languages (C/C++), require static program analysis and mostly lack generality.

Another approach uses something called needed-narrowing [3] or lazy instantiation [16] — where instances are only created when they are actually put to use by the program. This enables test data generators to adjust object instances during execution, when attempts are made to use them, so that they always have the required type. This technique is used by IRULAN [1] for generating tests in Haskell, which has lazy evaluation by default. For the purpose of complex type test data generation in Flycatcher, we will investigate whether we can mimic lazy evaluation for JavaScript objects in order to build test data with the appropriate complex type.

**Identifying bugs**

In dynamic languages such as JavaScript, the function signatures bear no type information and this makes it difficult to know whether an exception is raised due to a wrongly typed test argument or a true program bug. In the case where the exception is not a bug it could be due to two things: manipulating a badly initialised object or breaking a program precondition. The bad initialisation problem might be prevented through the use of lazy evaluation to initialise objects. The breaking of preconditions could be avoided by giving the tester the ability to impose restrictions on the test data generator, so that preconditions for the program are respected.

Being able to identify real program bugs during test data generation is an issue that will need further attention during the implementation of Flycatcher.

**Dealing with dynamically generated code**

Dynamic languages sometimes offer features that parse and evaluate a string at runtime and execute it as code, such as JavaScript's `eval` function. However, not only are these features potentially insecure, they make any analysis for test data generation much harder. As the general use of `eval` in JavaScript is prohibited anyway, we can safely ignore it for the purpose of Flycatcher.

## 2.3   Object-oriented test case generation

Most of the research on test generation focuses on testing imperative functions, such that the automated generation required is that of the functions' input parameters. However, when dealing with object-oriented code, a different approach is needed, as the unit under test changes from a function to an object. To test one of an object's methods three steps are necessary and should be repeated until the chosen coverage criterion for the method under test is met.

1. Instantiate the object

2. Call some of its methods to possibly modify its state

3. Assert that the method under test returns the expected answer

The method call sequences from the second step, as well as their inputs and the input to the constructor, should all evolve in a way that permits the successful coverage to be achieved. ATDG techniques that maximise coverage have been covered, but there has also been work done on how to guide the evolution of object-oriented tests cases *i.e.* including

## 2.4   JavaScript

### 2.4.1   Prominence

### 2.4.2   Idiosyncratic features

### 2.4.3   Object-oriented programming

# Chapter 3

# Project plan

# Chapter 4

# Evaluation strategy

# Bibliography

[1] ALLWOOD, T., CADAR, C., AND EISENBACH, S. High coverage testing of haskell programs. In *Proceedings of the 2011 International Symposium on Software Testing and Analysis* (2011), ACM, pp. 375–385.

[2] ANAND, S., PĂSĂREANU, C., AND VISSER, W. Jpf–se: A symbolic execution extension to java pathfinder. *Tools and Algorithms for the Construction and Analysis of Systems* (2007), 134–138.

[3] ANTOY, S., ECHAHED, R., AND HANUS, M. A needed narrowing strategy. In *Proceedings of the 21st ACM SIGPLAN-SIGACT symposium on Principles of programming languages* (1994), ACM, pp. 268–279.

[4] CADAR, C., DUNBAR, D., AND ENGLER, D. Klee: Unassisted and automatic generation of high-coverage tests for complex systems programs. In *Proceedings of the 8th USENIX conference on Operating systems design and implementation* (2008), USENIX Association, pp. 209–224.

[5] CADAR, C., GANESH, V., PAWLOWSKI, P., DILL, D., AND ENGLER, D. Exe: automatically generating inputs of death. *ACM Transactions on Information and System Security (TISSEC) 12*, 2 (2008), 10.

[6] DUCASSE, S., ORIOL, M., BERGEL, A., ET AL. Challenges to support automated random testing for dynamically typed languages.

[7] EDVARDSSON, J. A survey on automatic test data generation. In *Proceedings of the 2nd Conference on Computer Science and Engineering* (1999), no. x, pp. 21–28.

[8] GALLAGHER, M., AND LAKSHMI NARASIMHAN, V. Adtest: A test data generation suite for ada software systems. *Software Engineering, IEEE Transactions on 23*, 8 (1997), 473–484.

[9] GODEFROID, P., KLARLUND, N., AND SEN, K. Dart: directed automated random testing. In *ACM Sigplan Notices* (2005), vol. 40, ACM, pp. 213–223.

[10] GOLDBERG, A., WANG, T., AND ZIMMERMAN, D. Applications of feasible path analysis to program testing. In *Proceedings of the 1994 ACM*

*SIGSOFT international symposium on Software testing and analysis* (1994), ACM, pp. 80–94.

[11] GUPTA, N., MATHUR, A., AND SOFFA, M. Automated test data generation using an iterative relaxation method. In *ACM SIGSOFT Software Engineering Notes* (1998), vol. 23, ACM, pp. 231–244.

[12] HAN, S., AND KWON, Y. An empirical evaluation of test data generation techniques. *Journal of Computing Science and Engineering 2*, 3 (2008), 274–300.

[13] KING, J. A new approach to program testing. *Programming Methodology* (1975), 278–290.

[14] KING, J. Symbolic execution and program testing. *Communications of the ACM 19*, 7 (1976), 385–394.

[15] KOREL, B. Automated software test data generation. *Software Engineering, IEEE Transactions on 16*, 8 (1990), 870–879.

[16] LINDBLAD, F. Property directed generation of first-order test data. *TFP 7* (2007), 105–123.

[17] LLC, D. Programming language popularity, 2011.

[18] MAHMOOD, S. A systematic review of automated test data generation techniques. *School of Engineering, Blekinge Institute of Technology Box 520* (2007).

[19] MAIRHOFER, S. Search-based software testing and complex test data generation in a dynamic programming language. *Master's thesis, Blekinge Institute of Technology* (2008).

[20] MCMINN, P. Search-based software test data generation: a survey. *Software Testing, Verification and Reliability 14*, 2 (2004), 105–156.

[21] MEUDEC, C. Atgen: automatic test data generation using constraint logic programming and symbolic execution. *Software Testing, Verification and Reliability 11*, 2 (2001), 81–96.

[22] MICHAEL, C., AND MCGRAW, G. Automated software test data generation for complex programs. In *Automated Software Engineering, 1998. Proceedings. 13th IEEE International Conference on* (1998), IEEE, pp. 136–146.

[23] MICHAEL, C., MCGRAW, G., AND SCHATZ, M. Generating software test data by evolution. *Software Engineering, IEEE Transactions on 27*, 12 (2001), 1085–1110.

17

[24] MYERS, G., SANDLER, C., AND BADGETT, T. *The art of software testing.* Wiley, 2011.

[25] OFFUTT, A., JIN, Z., AND PAN, J. The dynamic domain reduction procedure for test data generation. *Software-Practice and Experience 29*, 2 (1999), 167–194.

[26] PĂSĂREANU, C., AND VISSER, W. A survey of new trends in symbolic execution for software testing and analysis. *International Journal on Software Tools for Technology Transfer (STTT) 11*, 4 (2009), 339–353.

[27] PLUQUET, F., MAROT, A., AND WUYTS, R. Fast type reconstruction for dynamically typed programming languages. In *ACM SIGPLAN Notices* (2009), vol. 44, ACM, pp. 69–78.

[28] PRASANNA, M., SIVANANDAM, S., VENKATESAN, R., AND SUNDARRAJAN, R. A survey on automatic test case generation. *Academic Open Internet Journal 15* (2005), 1–5.

[29] RAMAMOORTHY, C., HO, S., AND CHEN, W. On the automated generation of program test data. *Software Engineering, IEEE Transactions on*, 4 (1976), 293–300.

[30] SAI-NGERN, S., LURSINSAP, C., AND SOPHATSATHIT, P. An address mapping approach for test data generation of dynamic linked structures. *Information and Software Technology 47*, 3 (2005), 199–214.

[31] SEN, K., MARINOV, D., AND AGHA, G. *CUTE: A concolic unit testing engine for C*, vol. 30. ACM, 2005.

[32] TAHBILDAR, H., AND KALITA, B. Automated software test data generation: Direction of research. *International Journal of Computer Science and Engineering 2*.

[33] TILLMANN, N., AND DE HALLEUX, J. Pex–white box test generation for. net. *Tests and Proofs* (2008), 134–153.

[34] TRACEY, N., CLARK, J., AND MANDER, K. The way forward for unifying dynamic test-case generation: The optimisation-based approach. In *International Workshop on Dependable Computing and Its Applications (DCIA)* (1998), pp. 169–180.

[35] TRACEY, N., CLARK, J., MANDER, K., AND MCDERMID, J. An automated framework for structural test-data generation. In *Automated Software Engineering, 1998. Proceedings. 13th IEEE International Conference on* (1998), IEEE, pp. 285–288.

[36] VISVANATHAN, S., AND GUPTA, N. Generating test data for functions with pointer inputs. In *Automated Software Engineering, 2002. Proceedings. ASE 2002. 17th IEEE International Conference on* (2002), IEEE, pp. 149–160.

[37] ZHAO, R., AND LI, Q. Automatic test generation for dynamic data structures. In *Software Engineering Research, Management & Applications, 2007. SERA 2007. 5th ACIS International Conference on* (2007), IEEE, pp. 545–549.