



MASTER RESEARCH INTERNSHIP



INTERNSHIP REPORT

Adapting Amplified Unit Tests for Human Comprehension

Domain: Software Engineering — Artificial Intelligence

Author:
Simon BIHEL

Supervisor:
Benoit BAUDRY

KTH Royal Institute of Technology

Abstract:

TODO

Contents

Introduction	1
1 Background	1
1.1 Software Testing	1
1.1.1 Test Activities	1
1.1.2 Levels of Testing	2
1.1.3 Unit Testing	4
1.2 Elementary Metrics	6
1.3 Mutation Testing	7
1.4 The Need for Easy-to-use Tools	9
1.5 Cognitive Support Tool Development	9
2 Test Suite Amplification	10
2.1 Genetic Improvement	10
2.2 Test Amplification	10
2.3 DSpot	11
3 Problem Statement	12
3.1 The Need for Unit Test Cases Documentation	12
3.2 The Generated Random Noise	13
4 Related Works	13
4.1 Source Code Documentation	13
4.2 Source Code Change Documentation	14
4.3 Automatic Test Case Documentation	16
5 Contribution	18
5.1 Identifying Amplifications	18
5.2 Minimisation	18
5.3 Replace or Keep	18
5.4 Focus	18
5.5 Slicing	18
5.6 Natural Language Description	18
5.7 Ranking	19
6 Evaluation	19
6.1 Threat to Validity	19
6.2 Case Studies	19
6.3 Performances	19
Conclusion	19

Acknowledgments	19
References	19

Introduction

TODO

1 Background

In this section, we present the landscape this thesis fits into. We define testing (Section 1.1) as well as go over its use in the industry. In particular, to understand the needs of practitioners, we present how they assess the quality of their code (Sections 1.2 and 1.3) and what it takes for a new tool to be incorporated in their workflow (Sections 1.4 and ??).

1.1 Software Testing

In this section we give a definition for the test activities and their actors, the different abstraction levels of tests, and our precise object of study, unit tests.

TODO: say the Section 1.1.1 is a list of definitions to which the reader can refer to later on?

TODO: Why are we testing software and how do we do it

TODO: Say that the formal definitions aren't totally necessary

1.1.1 Test Activities

TODO: the text is too close to the oracle survey

Testing is about verifying that a system, for a given scenario, follows a certain behaviour that was previously defined. The *System-Under-Test* (SUT), in our domain a software system, has a set of components C . A scenario is sequence of stimuli that target a subset of C , and trigger responses: (i) feedbacks from the SUT; and (ii) changes in the state of its components. From [1], we use the following definition for to combination of stimuli and responses:

Definition 1.1 (Test Activities). For the SUT p , S is the set of stimuli that trigger or constrain p 's computation and R is the set of observable responses to a stimulus of p . S and R are disjoint. Test activities form the set $A = S \uplus R$. The disjoint union is used to label the elements of A .

The use of the terms “stimulus” and “observations” fits various test scenarios, functional and non-functional. As depicted in Figure 1, a stimulus can be either an explicit *input* from the tester, or an environmental factor that can affect the SUT. Stimuli include, among others, the platform, the language version, the resources available, interaction with an interface, etc. Observations encompass anything that can be discerned and measured like the content of a database, the visual responses on a computer screen, the time passed during the execution, etc.

Testing is about stimulation and observation, so we talk about *test activity sequences* that are comprised of at least one stimulus and one observation. But testing is also about verifying that the observed responses, the behaviour, match to ones that were previously defined. This checking part is done by another actor:

Definition 1.2 (Test Oracle). A test oracle $D : T_A \mapsto \mathbb{B}$ is a partial function from a test activity sequence to true or false.

An oracle can be defined using different methods: a set of specifications an expected value for a variable after the execution of a particular test activity sequence, the expectation of an absence of crash, etc.

A test oracle is a partial function because it can leave certain elements of the SUT's behaviour unspecified. This part of uncertainty is fundamental in the real world because systems have grown to be so complex that humans cannot anticipate all the reactions of the SUT in all environments. It is useful nonetheless to have a name for a theoretical oracle that would have an answer for every possible question:

Definition 1.3 (Ground Truth). The ground truth oracle, G , is a total test oracle that always gives an answer.

This allows us to define two properties for all oracles:

Definition 1.4 (Soundness). The test oracle D is sound for test activity a iff $D(a) \rightarrow G(a)$

Definition 1.5 (Completeness). The test oracle D is complete for test activity a iff $G(a) \rightarrow D(a)$

A test oracle cannot, in general, be complete as it would require handling every possible test activity sequence. It is not rare for an oracle to be unsound as humans can make errors writing specifications.

In practice, following pattern designs of modularity, the act of testing is composed of multiple test activity sequences that target specific elements of the SUT's behaviour:

Definition 1.6 (Test Case). A test case T_C is a test activities sequence that contains at least one stimulus and one observation with a test oracle that verifies each observation.

Definition 1.7 (Test Suite). A test suite is a collection of tests cases.

In other words [2], a test case is a subset of inputs, and a subset of specifications from the ground truth that are matched against the trace of the test's execution.

Another concept that we will use when talking about software evolution is regression testing [3]:

Definition 1.8 (Regression Testing). Regression testing is performed between two different versions of the SUT in order to provide confidence that newly introduced features do not interfere with the existing features.

It is a way to avoid having a human write formal specification but the problem remains of trying every possible scenario and observing any useful response.

TODO: give some clues to each concept to tell why it is a field of research?

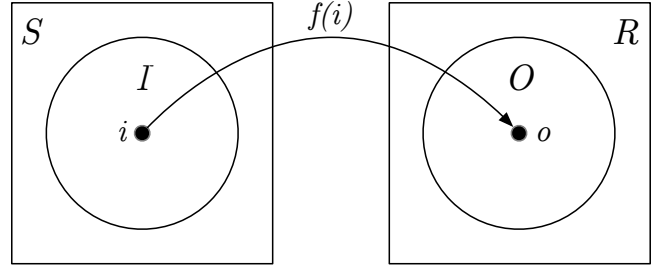


Figure 1: Stimuli and observations: S is anything that can change the observable behaviour of the SUT f ; R is anything that can be observed about the system's behaviour; I includes f 's explicit inputs; O is its explicit outputs; everything not in $S \cup R$ neither affects nor is affected by f .

1.1.2 Levels of Testing

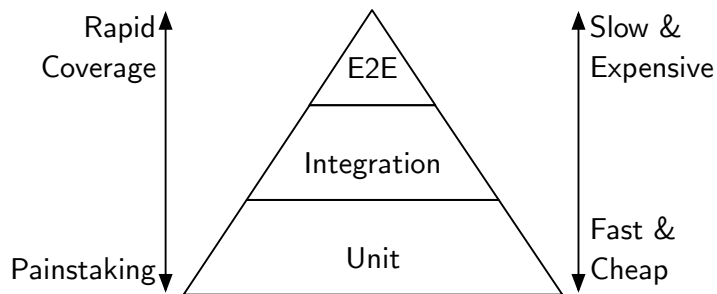


Figure 2: A view of the test pyramid.

In the case of software testing, the SUT can be different elements of a program, i.e. it can belong to different levels of abstraction. One can decide to test a whole program, by manually giving inputs and verifying that the output is correct. Or one could test individual functions, having a more control over the behaviour of each element.

These levels of test abstraction can be visualised with a test pyramid, as

shown in Figure 2. Tests are generally separated into three distinct categories. The top category is for the End-to-End tests, which means testing the program as a whole. It can mean interacting with the UI and evaluating the visual responses. Whilst it allows for directly evaluating the well behaviour of the program, each test is costly in time and resources because the whole program has to be run. Not only will each component be run, whether or not it requires testing, but time will also be spent on call to external libraries that are already tested. At the other end of the test spectrum are the unit tests. The goal is to test each and every component of the program (e.g. every function), as an individual and isolated from the rest of the program. Running a single component is faster than running the whole program, but require much more work from the person writing tests, as every component needs its set of tests. In between the two testing practices are the integration tests, which aim at ensuring the well collaboration of components *~work together as expected*, like verifying that gears rotate in the right way. Examples of integration tests include API testing, or simply the fact of calling a function using the result from another function call. This last example gives an intuition on how fuzzy the distinctions between all kinds of tests are, with definitions often relying on personal opinions¹². More details about the loose definition of unit tests are given Section 1.1.3.

The test hierarchy is often represented by a pyramid because it is generally advised to have many more unit tests than end-to-end tests. Practitioners often recommend putting an emphasis on unit test, with a distribution of 70% unit tests, 20% integration tests, and 10% end-to-end tests for example³. More details on the strengths and weaknesses of unit tests are given in the next section, but we can convey some general truths here. End-to-end tests require more maintenance as the software evolves, because test cases rely on many modules and their behaviour eventually change with time. Also, with a broad coverage, it is difficult to identify components at fault in the case of a failing test. Having unit tests for each component gives confidence that they behave accordingly and that if an integration test fails, the error lies in the interactions between the components. Again, the broader the test, the greater the number of interactions which makes it difficult identifying the one at fault. Because of this, the number of test levels has to increase alongside the software's growth — e.g. grouping components into modules, modules into services, and services as the system.

Programs can be tested manually by a human (e.g. by interacting with the UI), but tests are generally automated to save time. Automation often mean to programmatically describe a scenario, observations and their associated specifications. Different levels can require different tools and

¹<https://testing.googleblog.com/2009/07/software-testing-categorization.html>

²<https://martinfowler.com/bliki/IntegrationTest.html>

³<https://testing.googleblog.com/2015/04/just-say-no-to-more-end-to-end-tests.html>

frameworks⁴. For example, unit tests can be written in the same language as the main program, interacting with it like any other module (more details in the next section). But to test a graphical user interface, one might need to run a window system and control the cursor's movements and clicks.

1.1.3 Unit Testing

A unit test is supposed to target a precise component of the SUT, but defining the level of granularity is not trivial, even for experienced programmers⁵⁶ [4]. In the case of a program written with an Object-Oriented language, such as Java — we will put our focus on Java for the rest of this thesis — the targeted component can be a class or just a method of a class.

```
1 public class TreeListTest {  
2     public void testIterationOrder() {  
3         TreeList tl = new TreeList(10);  
4         for (int i = 0; i < 10; i++) {  
5             tl.add(i);  
6         }  
7         ListIterator it = tl.listIterator();  
  
8         int i = 0;  
9         while (it.hasNext()) {  
10            Integer val = it.next();  
11            assertEquals(i++, val.intValue());  
12        }  
13    }  
14 }
```

Listing 1: Example of an object-oriented unit test (taken, and adapted for readability, from the Apache Commons Collections, in the class `TreeListTest`, line 270⁷): it consists of test inputs (lines 3–7) that manipulate the SUT; and assertions (line 11). **TODO: what about lines 9 & 10**

TODO: the link might break in the future **TODO: kinda assumes that `.hasNext()` and `.next()`**

An example of unit test is given in Listing 1. Tests are usually sorted in classes, and each test case corresponds to a method (which is why test cases are sometimes referred to as test methods, as we will encounter later on). Specifications of expected values for observations are written using *assertions*. When the condition of an assertion is not met, we call it an assertion *failure*, it raises an exception[~] and causes the test case to fail. Assertions are a feature of test frameworks (i.e. internal DSL [5]), and allow test cases to be regular pieces of code. This has advantages and disadvantages, it is easy to write test cases but might not be the most friendly[~] way to express

⁴<https://martinfowler.com/articles/practical-test-pyramid.html>

⁵<https://martinfowler.com/bliki/TestPyramid.html>

⁶<https://martinfowler.com/bliki/UnitTest.html>

⁷<https://commons.apache.org/proper/commons-collections/xref-test/org/apache/commons/collections4/list/TreeListTest.html#L271>

a property on the software artifact targeted. Test methods often rely on naming conventions to help the reader understand what is the target, and what behaviour it is checking. Another pitfall to avoid is to test properties for the programming language used.

Unit tests are usually written by software engineers, as they are the one with the knowledge of the expected behaviour for the piece of software they have produced. Large companies will also have test engineers, focusing on broader tests and commercially important features, as well as software engineers working on test infrastructure and tooling⁸. Unit testing has becoming an integral part of modern software development with practices such as Test-driven Development (TDD)⁹ [6] that consists in writing test cases first, and then writing code that complies. Integrated development environments (IDEs) have been extended to make writing and reviewing tests easier. Test automation in general is also a pillar of modern software delivery cycles like DevOps¹⁰, with Continuous Integration (CI) tools that run test suites each time that code changes are pushed. Tests can then be refined using exploratory testing¹¹ [7], by trying to manually find flaws while documenting the process to then implement the results as test cases. Manual testing is particularly useful to rapidly detect issues like slow response time, misleading error messages and other design concerns.

Because unit test are written with the same programming language as their target, certain features of the language — that are supposed to help developers write clean code — can become a hindrance. A perfect example with Java are the private methods. They exist because they enforce the separation of the public interface and the inner workings of a class. But when one want to test a private method to make sure it does the right thing, it is not possible. The alternative is to extensively test the public methods, making sure they trigger all sorts of behaviour in the private part of the class. In this situation, we end testing a component that interacts with other components, as depicted in Figure 3, but then it becomes a small integration test¹². The granularity of tests does not only depend on the needs and desires of the developer, but also on the limitations of the language. Object-Oriented are particularly victim to this, as they are suited for composition design patterns [8], leaving no other option to the developer to try to limit interactions. Ways to avoid writing test that trigger a long chain of reactions is to use doubles (or mocks) which are hand written clones of classes instances. ~

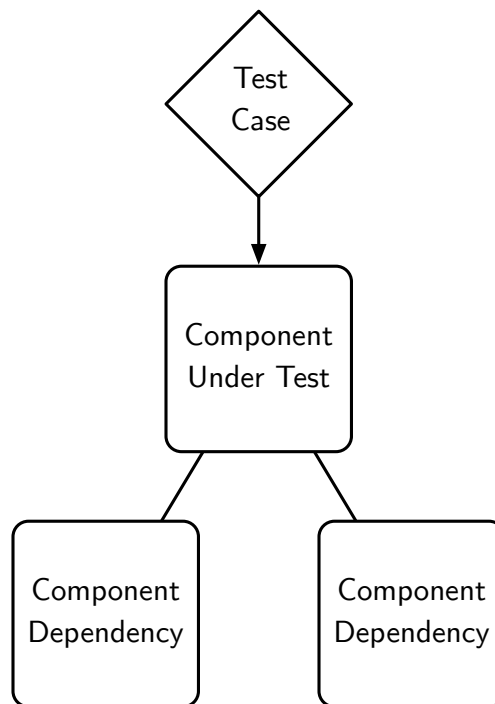


Figure 3: Test target’s dependencies.

Another common problem with testing are flaky tests. A flaky test is a test which result is not deterministic. There can be various causes for a flaky test, e.g. a direct usage of RNG function, or a crash happening only in scenario with a certain scheduling

⁸<https://testing.googleblog.com/2016/03/from-qa-to-engineering-productivity.html>

⁹https://en.wikipedia.org/wiki/Test-driven_development

¹⁰<https://en.wikipedia.org/wiki/DevOps>

¹¹<https://testing.googleblog.com/2008/05/exploratory-testing-on-chat.html>

¹²<https://martinfowler.com/bliki/UnitTest.html>

<code>public boolean add(String stringObject, Object value) {</code>	
<code> if (!notDeterministValues.contains(stringObject)) {</code>	1×
<code> if (observationValues.containsKey(stringObject)) {</code>	1×
<code> Object oldValue = observationValues.get(stringObject);</code>	1×
<code> if (oldValue == null) {</code>	1×
<code> if (value == null) {</code>	1×
<code> return true;</code>	1×
<code> } else {</code>	
<code> notDeterministValues.add(stringObject);</code>	!
<code> return false;</code>	!
<code> }</code>	
<code> } else if (!equals(oldValue, value)) {</code>	1×
<code> notDeterministValues.add(stringObject);</code>	1×
<code> return false;</code>	1×
<code> }</code>	
<code> } else {</code>	1×
<code> observationValues.put(stringObject, value);</code>	1×
<code> }</code>	
<code> return true;</code>	1×
<code> } else {</code>	
<code> return false;</code>	1×
<code> }</code>	
<code>}</code>	

Figure 4: Example of coverage reporting from the Coveralls¹⁵ tool. We can see that the statements in the first `else` branch have not been executed a single time

in the case of parallel applications. The advantage of small tests with controlled interaction with the environment is that they are less likely to be flaky. This a hard problem¹³ as flaky tests can be (i) hard to identify, (ii) hard to refactor, and (iii) hard to avoid in the first place.

TODO: how do you know you have a good oracle

TODO: junit ¹⁴

1.2 Elementary Metrics

Metrics have been developed to help developers automatically assess the quality of their code. Each metric measures one specific characteristic, such as the need for refactoring, the size of the program, or — and we are particularly interested in this one — the likelihood of bugs.

¹³<https://testing.googleblog.com/2017/04/where-do-our-flaky-tests-come-from.html>


The most basic metrics to measure how thoroughly tested a system is, are coverage based metrics. For example, during the test suite execution, one can keep track of all statements that were executed. At the end, you have a percentage of executed statements for the total number of statements. Instead of statements, one can also keep track of control flow branches explored. **TODO: give examples for other metrics?** These metrics, especially the statement coverage, are wide spread in the industry — with coverage tools integrated out-of-the-box with source code hosting services and continuous integration services. An example of coverage visual report is shown in Figure 4.

It is generally acknowledged that a system with high coverage means that the system is less likely to have bugs — but it is not foolproof [9, 10]. Such simple metrics are not good at pointing out corner cases. If we take the following example `if (C1 AND (C2 OR C3))` and `C3` makes the program crash, then it is possible to explore both branches without executing `C3`.

TODO: What are their limits

A drawback of metrics is that in practice, testers will change behaviour once they know how the measurement works¹⁶. They will write tests that will perform well for a given metric and ignoring other important factors (e.g. diversity¹⁷ [11, 12]). In these cases the metric has become the goal instead of a progress measurement.

1.3 Mutation Testing

Mutation testing is a method used to assess the quality of a test suite, giving it a grade called *mutation score*. The idea is to insert bugs in the software and see if these faulty versions still pass the tests or not. We call the derived versions *mutants*, hence the name mutation testing. When a mutant goes through the test suite without making a single test fail, it is considered *live*, and the opposite case it is considered *dead*. The global process is depicted in Figure 7. The main appeal is that, a test suite with a good mutation score, in addition to ensure the well behaviour of the SUT in useful scenarios, it can detect bad behaviours — which is particularly useful for regression testing.  There is a correlation between mutation score and coverage score as test suites with high mutation score tend to have good coverage [13]. It has been found that correlations between mutation scores and real fault detection are weak [14, 15].

TODO: a mutation covered line is stronger than a regular covered line

[foundational papers]

TODO: mutation testing has been extended in various directions, e.g. specification mutation, b

Mutants are obtained by applying *mutators* on the original version. These operators can be varied, from fine grained like changing a condition (e.g. substituting `≤` with `>`), to the granularity of classes [16] (e.g. delete the body of a method **[descartes]**). Running first mutation testing with large operators allows to rapidly detect gaps in the coverage of the test suite for unit testing but also for integration testing. Fine grained mutation can then be employed to meticulously evaluate the functional testing capabilities of the test suite [17]. Most of the time, we generate 1-*order* mutants, that means that a mutant is the result of a mutator applied once. A mutator will usually generate several mutants. For example, a mutator that replaces a `while` statement into a `do-while` statement will generate a mutant for each `while` in the artifact. In the case of *k*-order mutants,

¹⁵<https://coveralls.io/>

¹⁶<https://testing.googleblog.com/2009/09/7th-plague-and-beyond.html>

¹⁷<https://testing.googleblog.com/2009/06/by-james.html>

which can be thought of as the result of k successive 1-order mutants[18], the number of derived programs starts to blow up.

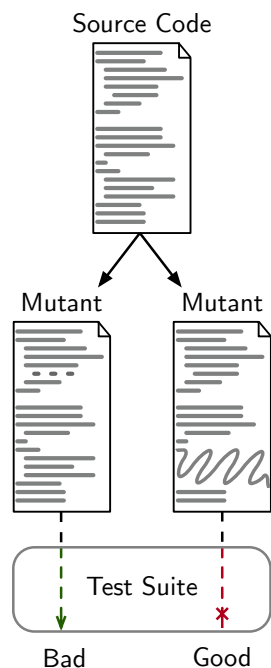


Figure 5: Mutation testing process.

Mutation testing suffers from drawbacks that have limited its use outside of the research world even though it has been an active research domain for many years [19]. The first obvious flaw is that it is slow. For compiled languages, all mutants have to be compiled, and that process often takes a lot of time for large programs. Then the whole test suite has to be executed for each mutant, which again is a long process. However, given certain trade-offs in terms of quality of the mutation analysis process (e.g. not generating all possible mutants), the computational complexity can be dodged [20, 21]. Optimized techniques for regression testing also exist [3], including minimisation, selection and prioritisation, which can reduce the execution time of the test suite. The mutant generation is also full of traps. The major pitfall that cannot be avoided is the equivalent mutant problem. Some mutants, although they have a different program than the original, have the same semantics. Which means that they will always pass the test suite and will not bring any insight on the SUT. The generation of mutants that are simply syntactically valid is also a non-trivial problem that requires pattern replacement instead of simple text replacement [22].

Another difficulty for mutation testing to take on in the industry is the lack of understanding by practitioners.~ The process of mutation and elimination can be confusing. But also, when a mutant is live, it is not always clear what actions should be taken in order to kill it.

TODO: maybe add more

TODO: teaching mutation testing using gamification

Multiple mutation tools exist, with differences on the level of abstraction they handle, different performances and different correlation with real fault detection. Many of them are robust thanks to the efforts from the research community to reach the industry and the need to have reliable implementations to do~ empirical studies related to mutation testing. One of these tools is PIT¹⁸ [23], a tool for the Java bytecode. The research version (i.e. state-of-the-art) of PIT is known to perform better than its competitors overall [24]. Two technical details are interesting for us. The first one is the way it stores the mutants generated. The location of a mutant is specified by (i) the name of the method and class, (ii) the method signature and (iii) the instruction on which the mutation occurs. This little information is sufficient to recreate each mutant while having a small storage footprint. The other important detail is that PIT is deterministic. That is, for the same program, the same mutants will always be generated, and thus the mutation score will always be the same.

The way mutation testing tool report their result is interesting. Instead of trying to shown each statement that has been mutated, and inspired by reports of line coverage, they usually display a regular line coverage report with shades of colour. The lightest means that the line has been executed, and the darker it gets, the higher the number of mutants located in this line. An example of mutation coverage visual report is shown in Figure 6.

¹⁸<http://pitest.org/>

¹⁹<http://code-defenders.org/>

```

55     public void call(int floor) {
56         if (floor >= 0 && floor <= topFloor) {
57             while (floor != currentFloor) {
58                 if (floor > currentFloor)
59                     goUp();
60                 else
61                     goDown();
62             }
63         }
64     }

```

Figure 6: Example of mutation coverage reporting from Code Defenders¹⁹ [25]. Shades of green are used to denote the amount of testing for a line.

1.4 The Need for Easy-to-use Tools

Software Engineering, as a research field, being tightly coupled with practices in the industry, it is important to understand the attempts at transferring research knowledge to the industry. There are many reasons to help the industry, for example: reducing the risks of bugs in critical systems, automating the development process, etc. One way researcher often go, is to develop a tool and then try to introduce it in the workflow of practitioners. But even with good intentions, the process of acceptance by practitioners for a research product is often complicated [26].

An issue that is often faced, is the difficulty for practitioners to understand the usage of the tool, or even its purpose. Thus, not only should a researcher consider the amount of bugs their tool can detect, they should also make sure it is accessible and easy to understand. Having an easy-to-use tool can also open new paths for research, revealing new problems, and making studies with professional developers easy.

Valuable experience has been earned from static analysis researchers’ attempts at integrating their tools in the workflow of their engineer colleagues [27, 28, 29]. Each false positive undermines the trust developers have in the tool. Each misunderstood error make developers angry. And each warning that do not *feel* relevant to the developer raise the likelihood that future warnings will be ignored. As humans, we avoid change unless it is necessary or is directly beneficial to us. Developers do not want to be slowed down by additional actors in their workflow, and companies do not want additional cost with little return.

The field of software testing, in which this thesis is set, is no stranger to this problem. Tests are, in most cases, written and maintained by general software engineers. Thus we must aim at producing tools with little overhead and fully integrated with popular test frameworks — making it easy for developers to try and see if the tool fits their need or not, and if they incorporate it in their workflow, it brings important data to demonstrate the scientific value of the tool.

1.5 Cognitive Support Tool Development

TODO Acceptance of a tool, outside of the research world, depends on many parameters, including its accessibility. In this section we describe how a tool can be designed to fit this requirement.

Research has been done on designing tool that do not overwhelm users (e.g. designing interfaces [30]) and theories on design processes have been developed. The grounded theory [31], for

example, argues for an iterative design process, by refining the product based on data and feedback from users.

Another method of evaluation that should not be dismissed are case studies [32]. **TODO**

2 Test Suite Amplification

In this Section, we provide more context for this thesis' contribution. We explain how global optimization techniques, including metaheuristic search techniques (e.g. genetic algorithms) can be applied to software engineering problems — for example, by using software metrics as target functions. Section 2.1 presents how source code can be manipulated for search processes, and Section 2.2 presents the specific problem of improving test suites. Section 2.3 introduces DSpot, a tool that enhances tests for a better coverage, and for which this thesis can be seen as an extension. ~

2.1 Genetic Improvement

TODO

Genetic improvement (GI) [33] uses automated search **TODO: explain what automated search is** to find improved versions of existing software. Applying search-based techniques for software engineering problems is not new [34] as is the field of genetic programming (GP) [35] (with, at the time of writing this thesis, is comprised of over 12000 references²⁰). But since ~ the name of GI has been coined in 2012 [36], the level of computing power reached has allowed research to achieve impressive results.

[foundational papers]

[the surprising creativity of digital evolution?]

Search operators are most frequently *delete*, *copy* and *replace* [37]. **TODO**

TODO: complex search techniques are not always the most efficient, and OO, again, pose certa

2.2 Test Amplification

As unit tests are program artefacts like others, GI can be applied on test suites [39] [A Systematic Literature Re and search techniques have been employed since 2005 [40]. It is relatively recent compared to works tackling the problem of Test Data Generation (i.e. generating tests from scratch for a piece of software), but it is only recent that we can expect software projects, large or modest, to come with strong test suites. Various goals can be pursued to improve a test suite, from automated refactoring to improve a specific metric, to enhancement of existing tests to fill the often partial coverage²¹. The term of test *amplification* is used to denote the process of processing and extending tests to reach further targets. Baudry et al. [40] wrote a tool for C# ~ to explore more of the input space and improve the mutation score. Yoo et al. [41] proposed the Test Data Regeneration (TDR) method for creating equivalent tests with different inputs — particularly useful for fault detection and avoiding overfitting **TODO: explain overfitting**. Xuan et al. [42, 43] created a tool to refactor test suite to improve dynamic analysis and fault detection.

²⁰<http://www.cs.bham.ac.uk/~wbl/biblio/>

²¹<https://testing.googleblog.com/2014/07/measuring-coverage-at-google.html>

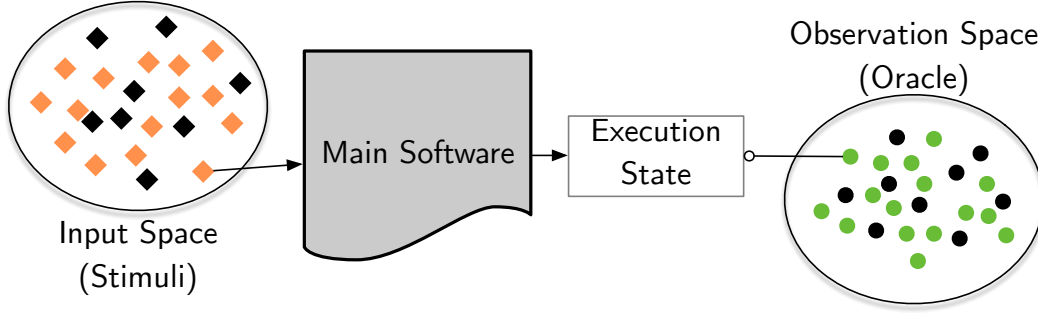


Figure 7: On the left, the testing input space is composed by specified input points (orange diamonds) and unspecified input points (black diamonds). On the right, the observation space over a given program state depicts the assertions of tests. The green circles are values that are already asserted in the existing test suite, the newly added assertions are shown as black circles.

2.3 DSpot

DSpot²² [Genetic-Improvement based Unit Test Amplification for Java] [44, 45, 46] is a test amplification tool. Its goal is to amplify tests individually to improve a given metric, e.g. mutation score. Going back to our definition of tests, with the notions of stimuli and observations, DSpot has a search process on both spaces. New inputs are generated, and new observations are made. A visual representation is given in Figure 7. Because of the oracle problem, DSpot focuses on regression testing — assertions (i.e. specifications) are generated from observations of an actual execution of the test. It is a Java tool and it uses Spoon [47] to manipulate the source code.

To explore more of the input space (we also talk about *amplifying* inputs), operators called I-Amplification are applied.

Amplification of literals Similarly to TDR, literals (numeric, boolean, string) are replaced by neighbours. For example, an integer can be multiplied by 2 or a random character can be added to a string.

Amplification of method calls Methods calls can be duplicated, removed, or made-up by picking a random variable in the test and calling one of its methods with random parameters.

Test objects If a new object is required for an amplified method call, one of the appropriate type will be built by using the default constructor.

Like 1-order mutants, only one amplification is applied at a time, to limit the number of generated tests. Because the random amplification can create flaky tests, they are executed three times to verify that they have a consistent behaviour.

As for the amplification of observations (A-Amplification), focus is put on checking more properties to kill more mutants. Assertions are generated as follows: the state of program is collected after execution of the test case to know the actual values and use these values as oracle.

Using mutation score as the target metric allows us to see the problem of amplification as a problem of finding bugs. In other words [48], the amplification process consists of improving the code coverage to execute faulty statements, and improving the propagation of faulty program states to an observable output.

DSpot is different from traditional test generator like EvoSuite [49] for multiple reasons. It is

²²<https://github.com/STAMP-project/dspot>

more than just using the existing test suite as a good starting population. **TODO** [50]

The result of DSpot is a set of amplified tests but it should not stop there. Building incrementally on top of humans work is valuable. Adding the generated tests in the main test suite has multiple advantages. First, is the review process. Developers have to review changes before it is accepted in the main code. During the review process developers can tweak the assertions, rearrange the inputs, etc. In that sense, DSpot fits as a software development bot, integrated in the workflow [51], that could be run on code changes, after making sure the tests suite is passing. The amplified tests would then be submitted as a Pull Request (or Merge Request) in the code version control platform used. Secondly, it is more efficient. **TODO**

3 Problem Statement

This thesis aim at helping developers understand amplified unit tests. There are two sides for this problem: supporting the developer understand the test (Section 3.1) and cleaning the tests from the noise injected during the amplification process (Section 3.2).

3.1 The Need for Unit Test Cases Documentation

The first challenge in rendering generated tests comprehensible, is to bundle them with explanations, like we do by documenting the code we write. Several works have already highlighted the need for test documentation, with developers surveys [52, 53, 54, 55, 56] and experiments [57] to measure the added value of documentation. Even when they are short, tests can be complex and hard to understand. In all cases, having documentation — which include documentation comments like JavaDoc, but also the naming of the tests and variable — has many advantages [58] **[How Do Automatically Generated Unit Tests Influence Software Maintenance?]** :

- faster to get familiar with the test;
- faster fault localisation (e.g. figuring out if there is a bug in the test or the SUT) when the test fails; and
- helps building trust in a test generator if it can offer a proof for its result.

Information to understand can be divided into four categories: (i) the component targeted; (ii) the context of the test method; (iii) the actual piece of code (i.e. the logic) that is the test method; and (iv) the properties checked. And we can split them again in two groups: the *what* and *why* informations. As we will see in Section 4, extensive works have been produced to generate automatic documentation of the what information. As for the why information, it is difficult to automatically extract the motivation of the developer that wrote the test. But in the case of generated test, we have the reason why they were kept: they bring something to the metric. Thus, if we focus on amplification for mutation score, we would need to explain this score and what new mutants are killed. A lack of works on assisting humans to get the grasp of a mutation score makes a good case for us to tackle the problem of explain a killed mutant.

In source code management platforms, pull requests (through which amplified are submitted for review and merge in the main branch) are accompanied by textual descriptions. Given that amplified tests should, in most cases, extend an existing test — that we can suppose have documentation already — the explanation of the added value of the amplification would fit in the pull request description. It is important to restrict the scope to simple text, as it would not require a change in developers workflow.

Structural		Collaborational	Creational
Accessor	Mutator		
Get	Set	Collaborator-Accessor	Constructor
Predicate	Command	Collaborator-Mutator	Copy-Constructor
Property			Destructor
			Factory

Table 1: A taxonomy of method stereotypes.

3.2 The Generated Random Noise

Through the search process, noise is injected in the test case. This problem of ugly generated tests code is a shared problem between all kinds of generators [59]. Because documentation is not sufficient to understand and review the code, it needs to be clean and logical to help the comprehension.

On a side note, in many modest projects, unit test are used as examples of the program’s usage, and sometimes they replace traditional documentation. In those cases, developers put a lot of efforts in making tests easy to understand. Doing part of that chore automatically would remove part of the burden for humans.

4 Related Works

In this section we review the large body of related works from the domain of software maintenance [60]. We are interested in tools that automatically generate documentation — and there are many of them, as it is well known that software documentation is important [61] (even for modern agile methods) but it is time-consuming, often times poor and incomplete [62], and is hard to maintain [63]. They range from documenting general purpose source code (Section 4.1), to commit message generation (Section 4.2), and documenting test cases (Section 4.3).

4.1 Source Code Documentation

TODO: these works are useless, why describe them?

Works on source code documentation — also called software artifact summarization — often consist of retrieving information from the source code to then arrange this information in a way that is easier and quicker to grasp for a human. This section does not aim at providing a survey of the field (for which there are already, including automatic summarization [64, 65], natural language paraphrasing [66] and automatic comment generation [67]), but to present the general techniques used as well as the kind of results those tools produce.

As said above, such tools will rearrange information, stripping irrelevant pieces, summarizing chunks on code, and put central information on the front. In other words, they render the what information in a human comprehensible way — they *paraphrase* the code. This is the main reason why, even automatic summarization of natural language documents has been attempted by researchers for more than half a century [64], finding a link with our problem of explaining the why information is not trivial.

One of the most influential works was produced by Dragan et al. [68, 69], which is about method stereotypes. They proposed a taxonomy of C++ methods, given in Table 1, and a way to automatically assign stereotypes. This way of annotating methods lead to other works — using stereotype distribution as a signature descriptor for systems [70], extending the taxonomy for Java methods and classes [71].

Buse and Weimer [72] wrote a tool to automatically document the exceptions a method can raise. Exceptions are dugged up from the call stack and the conditions for them to be raised are translated in natural language. They compared their generated documentation with hand-written documentation in a survey with developers.


Hill [73, 74] has developed a model, the Software Word Usage Model (SWUM), that captures the occurrences of words in code, and also their linguistic and structural relationships. “SWUM captures the conceptual knowledge of a programmer expressed through both natural language information and programming language structure and semantics, in the form of phrasal concepts. In contrast to existing lexical approaches, we take a transformative step in automatic analysis for software engineering tools by accounting for how words occur together in code, rather than just counting their frequencies.” Because programmers tend to choose descriptive and meaning names for code units [75], natural language descriptions of code artefacts can be made, without necessarily having to understand the code.

Sridhara et al. [76, 77] worked on automatically generating summary comments for Java methods. For that they wrote an algorithm to extract important code statements and used SWUM to paraphrase those statements. They did a survey with developers to evaluate the accuracy and usefulness of their contribution.


McBurney and McMillan [78] wrote a tool to generate natural language descriptions of the context of Java methods. What they call context is how the methods are invoked. They use the SWUM to describe the method, describe the method’s caller, what the return value is used for, what methods it calls, and how the return value can be used. They did two user studies.

Sridhara et al. [79] worked on describing methods from a high level point of view. They detect high level actions (i.e. statements grouped by a loop or a condition) and use the SWUM to describe them. They did a survey with developers to judge the descriptions synthesised.

Moreno et al. [80] developed a tool to automatically generate summaries of Java classes. A summary is composed of the class’ stereotype [71], a description of its interface, structure and behaviour (by enumerating the most important methods). Summaries are in natural language and generated using templates. They did a survey with developers to judge how readable and understandable the summaries are.

Rastkar et al. [81] worked on summarizing bug reports, including involved discussions. They label sentences and extract individual information. Works related to bug reports are important for us as a killed mutant is a bug in a certain way — but the knowledge we can extract from these works is not obvious because part of the job of a bug report is to help localise the bug, and in the case of mutation testing we already know the source. Nevertheless, we can learn from these works how to explain the repercussions of a mutant. Except that in this case, it is focused on summarizing a discussion. 

4.2 Source Code Change Documentation

Software evolves, iteratively, line by line, commit by commit.  Past the phase of assimilation when a developer first start to work a on project, a significant part of the job consists of managing

Categories of commits	Description
Implementation	New requirements
Corrective	Processing failure; performance failure; implementation failure
Adaptive	Change in data environment
Perfective	Processing inefficiency; performance enhancement; maintainability
Non-functional	Code clean-up; legal; source control system management

Table 2: Categories of commits in terms of maintenance task and corresponding description.

the growth of the project, making small changes to the code base without unwanted side effects, and reviewing commits. Commits bundle a cohesive set of changes with a textual description, in which it is customary to provide some context, the motivation and possibly a digest of the changes — part what and part why information. Although developers view accurate and concise descriptions as useful, few put the efforts in writing such comments. In consequence, researchers have written tools to automatically generate commit messages.

Dragan et al. [82] proposed commit stereotypes. Those stereotypes are composed of the modified classes and methods’ stereotypes. For example, a commit where 75% of the added/deleted methods are factory methods denotes an “object creation modifier”.

Cortés-Coy et al. [83, 84] developed *ChangeScribe*, a commit message generator. From Git²³ changes they identify the responsibilities of the modified classes and methods, attribute stereotypes to changes [82], keep the most impactful (i.e. the most called modified methods or classes) changes, and generate natural language descriptions. Messages are generated using templates. The evaluated their tool with survey, comparing real-life commit messages and generated messages, asking the developers which one they preferred. They also provide a table listing many approaches for generating descriptions of source code changes and software artifacts, providing more references than this section.

Jiang and McMillan [85] worked on generating one-phrase summaries of commits. They used a classifier trained on real commits to label a commit with a verb (e.g. add, fix), identify the most important words occurring in the `diff`, and then use common phrase structure to combine the verb and the subject.

Jiang et al. [86] used neural machine translation to generate a commit message from a `diff`. Neural networks were also used by Loyola et al. [87].

Shen et al. [88] worked on summarization of what and why information. The what information is generated by paraphrasing the `diff`, dropping some information or adding some depending of how important the modified unit is. The why information is generated using categories proposed by Swanson [60], given in Table 2, and classify commits using the rules proposed by Moreno et al. [71]. The sentences were generated using templates.

Buse and Weimer [89] produced a technique for documenting program changes, using symbolic execution. They use a minimalist programmatic syntax to express changes (e.g. change in condition) and paraphrase changes at the statement level, instead of line granularity. Messages are generated as a pattern-matching and replace. They also studied some open-source projects and showed that what information 12% more common than why information in commit messages. They conducted a human study to compare their tool against existing human-written messages.

²³<https://git-scm.com/>

Buckley et al. [90] also proposed a taxonomy of software changes. Their taxonomy focuses exclusively on the what information of a change, the mechanisms of change and the factors that influence these mechanisms. Mechanisms include, for example, the change type or the change history, while the influencing factors include, for example, the change frequency, the granularity or the change propagation. In other words, they focus on the *when*, *where*, *what*, and *how* aspects.

Rastkar et al. [91] produced a generator of feature summary for Java projects, to help programmers be aware of, and understand, the ramifications of a feature. A summary provides the number of methods implementing the feature, objects supported, dependencies, and names of methods involved. Relationships between code elements are extracted, including calls between methods, field declarations, class extensions, etc. NL sentences are generated with templates. They evaluated their tool with a lab study, asking developers to modify or add a feature using the summaries, observing their usage and surveyed the developers at the end.

4.3 Automatic Test Case Documentation

Finally, the piece of work that is most related to this thesis is about documentation for test cases. Because of test cases are one precise kind of software artefact, the works that we have covered in the previous two sections can be cumbersome here (e.g. stereotypes are not suited).

As unit tests validate a specific component, in object-oriented languages like Java, they often consist in calling a few methods, multiple times, that will affect the state of an object and then verify that the object is in the right state. This kind of method is called a *focal* method, and Ghafari et al. [92] have developed a way to automatically identify them. They use static data-flow analysis to examine the call graph and extract the last method that has a side-effect on the assertion. This work is particularly useful to summarise test cases, helping to identify the key stimulus, but it is not directly useful for us as we want to explain a corner case (i.e. an amplification that kills a new mutant).~

Cornelissen et al. [93] propose the use of sequence diagrams²⁴ to visualize Java test cases. Because test suites are sometimes the entry points for new contributors to understand the scope and the inner workings of a project, visualizing the actors and their interactions can be more useful than a simple class diagram. Diagrams can be used for complex test cases as they hide certain call (e.g. constructors) and have limited call stack depth.

According to Meszaros [94], a unit test case typically has four phases: setup, exercise, verify and teardown. The setup phase instantiates the classes that the test uses. The exercise phase is the stimulation. The verify phase is the use of the oracle with assertions. The teardown phase is for the cleanup.~ **TODO: move that elsewhere**

Kamimura et al. [95] have developed a generator of textual summaries for Java test cases. A summary is composed of NL paraphrases of the most unique method invocations and the assertions. There are some interesting details, for example, when a variable's name is too short (e.g. "a1") they make-up a new name with the variable's type. Sentences are built using pre-defined templates.

Li et al. [56] have written a C# tool that generates descriptions for test methods. The output, a GUI window, is comprised of 4 pieces of information:

1. a general sentence describing the purpose of the test method (based on class, method and arguments signatures using the SWUM approach [96]);
2. descriptions of the focal methods (i.e. the line number and the assertion it is linked to);

²⁴https://en.wikipedia.org/wiki/Sequence_diagram

3. descriptions of the assertions (i.e. natural language (NL) paraphrasing); and
4. the slicing paths (backward slicing [97]) for the variables validated with an assertions (i.e. NL paraphrasing of each step).

NL sentences are generated with templates.

Li [98] also proposed 21 stereotypes for methods in Java unit tests, based on the jUnit²⁵ API (14 stereotypes) but also using the data and control flow (7 stereotypes). The API-based stereotypes allow the extraction of fine grained what information — for example, with `fail` (which raises a failure if the test cases reaches that statement) you can deduce that the test is a “utility verifier” or that a method with the annotation `@Before` is a “test initializer”. On the other hand, data/control flow analyses can capture the context of the API calls. The purpose of the test can be deduced, for example, if the number of assertions within branch conditions is greater than 0, it means that the test is a “branch verifier”. It is also possible to get insights in the context of the test case in relation to the test class, for example, a test case is a “public field verifier” if values in assertions are from public field accesses.

Beck et al. [99] wrote a framework that generates method execution reports, mostly textual but with a few embedded visualisations. The report includes a summary (number of calls, recursion), summaries for each method calls (number of calls, recursion depth), time consumption tear down, and the code. Embedded visualisations consists of fills bars for relative quantities and histograms for recursion stack. The text generation involves sentence templates and decision graphs (e.g. if there is recursion then present recursion information). All data are collected dynamically.

Delamaro and Maldonado [100] developed a fully fledged GUI for mutation testing, providing test cases execution and management, mutants creation and management, and mutation score reports. A report includes general information like the test cases and the killed mutant’s names, the compilation command used, etc.

CODE DEFENDERS²⁶ [25] is a game²⁷ that aims at teaching [101] software testing concepts, especially mutation testing. In this game, players are either attackers or defenders. Attackers write mutants, and defenders write tests cases. If there are more killed mutants than live ones the defenders are the winner, and conversely if there are more live mutants. The way a mutant is presented is with a simple `diff` (i.e. the modified statement is highlighted). This representation makes sense in a GUI with mutants on one side and the tests on the other side. Another motivation is that the game aims at teaching the current state of mutation testing. Because of these two choices, it is hard for us to draw inspiration.

Qusef et al. [102] worked on keeping track of traceability between Java unit tests and tested class, to maintain consistency in test classes and give some hints on the target of a test case. They use dynamic slicing [103] to identify all the classes that affect the result of the last assert statement in each unit-test method. And, because the slicing is an overestimate, they use conceptual coupling [104] (analysing names and comments) to filter out superfluous classes (e.g. mock objects). This problem of traceability is useful for us as we would like the amplified tests to remain in the same scope as the original, aiming for a more through coverage of the target instead of just random mutants.

Most of the papers related to human assistance reviewed here were evaluated with developers surveys. Proteum/IM 2.0 [100] used a case study, and CODE DEFENDERS conducted a lab study.

²⁵<https://junit.org/>

²⁶<http://code-defenders.org/>

²⁷https://en.wikipedia.org/wiki/Gamification_of_learning

5 Contribution

TODO

5.1 Identifying Amplifications

TODO

5.2 Minimisation

TODO TODO: put slicing before?

TODO: removing useless assertions

TODO: cannot use general purpose techniques[105, 106] because we want the original part intact

5.3 Replace or Keep

TODO

5.4 Focus

TODO [107]

5.5 Slicing

TODO²⁸ [108]

TODO: srcSlice²⁹ [109] does not support Java (yet)

5.6 Natural Language Description

TODO

TODO: Focus on mutation testing

TODO: avoid talking about mutants TODO: How to describe a mutant
[110, 111, 112]

TODO Elements at our disposal:

1. the original test;
2. the amplified test (its AST representation, with annotated amplified nodes); and
3. the newly killed mutants.

List of possible information:

- control flow branches
- name of mutant
- code under test [102]
- traceability link
- label/stereotype

TODO: on the usefulness of nl

TODO: what can we draw from the related works TODO: assuming knowledge of original test

²⁸<http://wala.sourceforge.net/>

5.7 Ranking

TODO

6 Evaluation

TODO

6.1 Threat to Validity

TODO

TODO: DSpot is not yet established and recognized in the community. It is difficult to have

6.2 Case Studies

TODO

6.3 Performances

TODO

Conclusion

TODO

Acknowledgments

Thanks to Benoit Baudry and Martin Monperrus for their guidance. Thanks to Benjamin Danglot for his collaboration and all his work on DSpot. Thanks to Zimin Chen, Nicolas Harrand, He Ye and Long Zhang for making daily life enjoyable. This internship was supported by the Fondation Rennes 1 and its patrons. Many thanks to KTH for hosting me.

References

- [1] E. T. Barr, M. Harman, P. McMinn, M. Shahbaz, and S. Yoo, “The oracle problem in software testing: A survey,” *IEEE transactions on software engineering*, vol. 41, no. 5, pp. 507–525, 2015.
- [2] G. Bernot, M. C. Gaudel, and B. Marre, “Software testing based on formal specifications: a theory and a tool,” *Software Engineering Journal*, vol. 6, no. 6, pp. 387–405, 1991.
- [3] S. Yoo and M. Harman, “Regression testing minimization, selection and prioritization: a survey,” *Software Testing, Verification and Reliability*, vol. 22, no. 2, pp. 67–120, 2012.
- [4] P. Runeson, “A survey of unit testing practices,” *IEEE software*, vol. 23, no. 4, pp. 22–29, 2006.

- [5] M. Fowler, *Domain-specific languages*. Pearson Education, 2010.
- [6] K. Beck, *Test-driven development: by example*. Addison-Wesley Professional, 2003.
- [7] C. Kaner, J. Falk, and H. Q. Nguyen, *Testing Computer Software Second Edition*. Dreamtech Press, 2000.
- [8] P. Wolfgang, “Design patterns for object-oriented software development,” *Reading Mass*, p. 15, 1994.
- [9] D. Hovemeyer and W. Pugh, “Finding bugs is easy,” *Acm sigplan notices*, vol. 39, no. 12, pp. 92–106, 2004.
- [10] L. Inozemtseva and R. Holmes, “Coverage is not strongly correlated with test suite effectiveness,” in *Proceedings of the 36th International Conference on Software Engineering*, pp. 435–445, ACM, 2014.
- [11] B. Baudry, M. Monperrus, C. Mony, F. Chauvel, F. Fleurey, and S. Clarke, “Diversify: Ecology-inspired software evolution for diversity emergence,” in *Software Maintenance, Reengineering and Reverse Engineering (CSMR-WCRE), 2014 Software Evolution Week-IEEE Conference on*, pp. 395–398, IEEE, 2014.
- [12] B. Baudry and M. Monperrus, “The multiple facets of software diversity: Recent developments in year 2000 and beyond,” *ACM Computing Surveys (CSUR)*, vol. 48, no. 1, p. 16, 2015.
- [13] B. Assylbekov, E. Gaspar, N. Uddin, and P. Egan, “Investigating the correlation between mutation score and coverage score,” in *Computer Modelling and Simulation (UKSim), 2013 UKSim 15th International Conference on*, pp. 347–352, IEEE, 2013.
- [14] M. Papadakis, D. Shin, S. Yoo, and D.-H. Bae, “Are mutation scores correlated with real fault detection? a large scale empirical study on the relationship between mutants and real faults,” in *40th International Conference on Software Engineering, May 27-3 June 2018, Gothenburg, Sweden*, 2018.
- [15] R. Just, D. Jalali, L. Inozemtseva, M. D. Ernst, R. Holmes, and G. Fraser, “Are mutants a valid substitute for real faults in software testing?,” in *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering*, pp. 654–665, ACM, 2014.
- [16] S. Segura, R. M. Hierons, D. Benavides, and A. Ruiz-Cortés, “Mutation testing on an object-oriented framework: An experience report,” *Information and Software Technology*, vol. 53, no. 10, pp. 1124–1136, 2011.
- [17] W. E. Howden, “Weak mutation testing and completeness of test sets,” *IEEE Transactions on Software Engineering*, no. 4, pp. 371–379, 1982.
- [18] K. Wah, “A theoretical study of fault coupling,” *Software testing, verification and reliability*, vol. 10, no. 1, pp. 3–45, 2000.

- [19] Y. Jia and M. Harman, “An analysis and survey of the development of mutation testing,” *IEEE transactions on software engineering*, vol. 37, no. 5, pp. 649–678, 2011.
- [20] A. J. Offutt, G. Rothermel, and C. Zapf, “An experimental evaluation of selective mutation,” in *Proceedings of the 15th international conference on Software Engineering*, pp. 100–107, IEEE Computer Society Press, 1993.
- [21] J. Mořucha and B. Rossi, “Is mutation testing ready to be adopted industry-wide?,” in *International Conference on Product-Focused Software Process Improvement*, pp. 217–232, Springer, 2016.
- [22] A. Simão, J. C. Maldonado, and R. da Silva Bigonha, “A transformational language for mutant description,” *Computer Languages, Systems & Structures*, vol. 35, no. 3, pp. 322–339, 2009.
- [23] H. Coles, T. Laurent, C. Henard, M. Papadakis, and A. Ventresque, “Pit: a practical mutation testing tool for java,” in *Proceedings of the 25th International Symposium on Software Testing and Analysis*, pp. 449–452, ACM, 2016.
- [24] M. Kintis, M. Papadakis, A. Papadopoulos, E. Valvis, N. Malevris, and Y. Le Traon, “How effective are mutation testing tools? an empirical analysis of java mutation testing tools with manual analysis and real faults,” *Empirical Software Engineering*, pp. 1–38, 2017.
- [25] J. M. Rojas, T. D. White, B. S. Clegg, and G. Fraser, “Code defenders: Crowdsourcing effective tests and subtle mutants with a mutation testing game,” in *Proceedings of the 39th International Conference on Software Engineering*, pp. 677–688, IEEE Press, 2017.
- [26] M. Delahaye and L. Bousquet, “Selecting a software engineering tool: lessons learnt from mutation analysis,” *Software: Practice and Experience*, vol. 45, no. 7, pp. 875–891, 2015.
- [27] A. Bessey, K. Block, B. Chelf, A. Chou, B. Fulton, S. Hallem, C. Henri-Gros, A. Kamsky, S. McPeak, and D. Engler, “A few billion lines of code later: using static analysis to find bugs in the real world,” *Communications of the ACM*, vol. 53, no. 2, pp. 66–75, 2010.
- [28] C. Sadowski, J. Van Gogh, C. Jaspan, E. Söderberg, and C. Winter, “Tricorder: Building a program analysis ecosystem,” in *Proceedings of the 37th International Conference on Software Engineering-Volume 1*, pp. 598–608, IEEE Press, 2015.
- [29] C. Sadowski, E. Aftandilian, A. Eagle, L. Miller-Cushon, and C. Jaspan, “Lessons from building static analysis tools at google,” *Communications of the ACM*, vol. 61, no. 4, pp. 58–66, 2018.
- [30] S. Oviatt, “Human-centered design meets cognitive load theory: designing interfaces that help people think,” in *Proceedings of the 14th ACM international conference on Multimedia*, pp. 871–880, ACM, 2006.
- [31] K.-J. Stol, P. Ralph, and B. Fitzgerald, “Grounded theory in software engineering research: a critical review and guidelines,” in *Software Engineering (ICSE), 2016 IEEE/ACM 38th International Conference on*, pp. 120–131, IEEE, 2016.

- [32] B. Flyvbjerg, “Five misunderstandings about case-study research,” *Qualitative inquiry*, vol. 12, no. 2, pp. 219–245, 2006.
- [33] J. Petke, S. Haraldsson, M. Harman, D. White, J. Woodward, *et al.*, “Genetic improvement of software: a comprehensive survey,” *IEEE Transactions on Evolutionary Computation*, 2017.
- [34] P. McMinn, “Search-based software testing: Past, present and future,” in *Software testing, verification and validation workshops (icstw), 2011 ieee fourth international conference on*, pp. 153–163, IEEE, 2011.
- [35] J. R. Koza, “Genetic programming as a means for programming computers by natural selection,” *Statistics and computing*, vol. 4, no. 2, pp. 87–112, 1994.
- [36] M. Harman, W. B. Langdon, Y. Jia, D. R. White, A. Arcuri, and J. A. Clark, “The gismoe challenge: Constructing the pareto program surface using genetic programming to find better programs (keynote paper),” in *Automated Software Engineering (ASE), 2012 Proceedings of the 27th IEEE/ACM International Conference on*, pp. 1–14, IEEE, 2012.
- [37] J. Petke, “New operators for non-functional genetic improvement,” in *Proceedings of the Genetic and Evolutionary Computation Conference Companion*, pp. 1541–1542, ACM, 2017.
- [38] S. Shamshiri, J. M. Rojas, L. Gazzola, G. Fraser, P. McMinn, L. Mariani, and A. Arcuri, “Random or evolutionary search for object-oriented test suite generation?,” *Software Testing, Verification and Reliability*, 2017.
- [39] B. Danglot, O. Vera-Perez, Z. Yu, M. Monperrus, and B. Baudry, “The emerging field of test amplification: A survey,” *arXiv preprint arXiv:1705.10692*, 2017.
- [40] B. Baudry, F. Fleurey, J.-M. Jézéquel, and Y. Le Traon, “Automatic test case optimization: A bacteriologic algorithm,” *ieee Software*, vol. 22, no. 2, pp. 76–82, 2005.
- [41] S. Yoo and M. Harman, “Test data regeneration: generating new test data from existing test data,” *Software Testing, Verification and Reliability*, vol. 22, no. 3, pp. 171–201, 2012.
- [42] J. Xuan, B. Cornu, M. Martinez, B. Baudry, L. Seinturier, and M. Monperrus, “Dynamic analysis can be improved with automatic test suite refactoring,” *arXiv preprint arXiv:1506.01883*, 2015.
- [43] J. Xuan, B. Cornu, M. Martinez, B. Baudry, L. Seinturier, and M. Monperrus, “B-refactoring: Automatic test code refactoring to improve dynamic analysis,” *Information and Software Technology*, vol. 76, pp. 65–80, 2016.
- [44] B. Baudry, S. Allier, M. Rodriguez-Cancio, and M. Monperrus, “Automatic software diversity in the light of test suites,” *arXiv preprint arXiv:1509.00144*, 2015.
- [45] B. Baudry, S. Allier, and M. Monperrus, “Tailored source code transformations to synthesize computationally diverse program variants,” in *Proceedings of the 2014 International Symposium on Software Testing and Analysis*, pp. 149–159, ACM, 2014.
- [46] B. Baudry, S. Allier, M. Rodriguez-Cancio, and M. Monperrus, “Dspot: Test amplification for automatic assessment of computational diversity,” *arXiv preprint arXiv:1503.05807*, 2015.

- [47] R. Pawlak, M. Monperrus, N. Petitprez, C. Noguera, and L. Seinturier, “Spoon: A Library for Implementing Analyses and Transformations of Java Source Code,” *Software: Practice and Experience*, vol. 46, pp. 1155–1179, 2015.
- [48] S. Shamshiri, R. Just, J. M. Rojas, G. Fraser, P. McMinn, and A. Arcuri, “Do automatically generated unit tests find real faults? an empirical study of effectiveness and challenges (t),” in *Automated software engineering (ASE), 2015 30th IEEE/ACM international conference on*, pp. 201–211, IEEE, 2015.
- [49] G. Fraser and A. Arcuri, “Evosuite: automatic test suite generation for object-oriented software,” in *Proceedings of the 19th ACM SIGSOFT symposium and the 13th European conference on Foundations of software engineering*, pp. 416–419, ACM, 2011.
- [50] J. M. Rojas, M. Vivanti, A. Arcuri, and G. Fraser, “A detailed investigation of the effectiveness of whole test suite generation,” *Empirical Software Engineering*, vol. 22, no. 2, pp. 852–893, 2017.
- [51] S. Urli, Z. Yu, L. Seinturier, and M. Monperrus, “How to design a program repair bot? insights from the repairnator project,” in *ICSE 2018-40th International Conference on Software Engineering, Track Software Engineering in Practice (SEIP)*, pp. 1–10, 2018.
- [52] E. Daka and G. Fraser, “A survey on unit testing practices and problems,” in *Software Reliability Engineering (ISSRE), 2014 IEEE 25th International Symposium on*, pp. 201–211, IEEE, 2014.
- [53] M. P. Prado, E. Verbeek, M.-A. Storey, and A. M. Vincenzi, “Wap: Cognitive aspects in unit testing: The hunting game and the hunter’s perspective,” in *Software Reliability Engineering (ISSRE), 2015 IEEE 26th International Symposium on*, pp. 387–392, IEEE, 2015.
- [54] M. P. Prado and A. M. Vincenzi, “Advances in the characterization of cognitive support for unit testing: The bug-hunting game and the visualization arsenal,” in *Software Reliability Engineering Workshops (ISSREW), 2016 IEEE International Symposium on*, pp. 213–220, IEEE, 2016.
- [55] M. P. Prado and A. M. R. Vincenzi, “Towards cognitive support for unit testing: a qualitative study with practitioners,” *Journal of Systems and Software*, 2018.
- [56] B. Li, C. Vendome, M. Linares-Vásquez, D. Poshyanyk, and N. A. Kraft, “Automatically documenting unit test cases,” in *Software Testing, Verification and Validation (ICST), 2016 IEEE International Conference on*, pp. 341–352, IEEE, 2016.
- [57] S. Panichella, A. Panichella, M. Beller, A. Zaidman, and H. C. Gall, “The impact of test case summaries on bug fixing performance: An empirical investigation,” in *Software Engineering (ICSE), 2016 IEEE/ACM 38th International Conference on*, pp. 547–558, IEEE, 2016.
- [58] E. Daka, J. M. Rojas, and G. Fraser, “Generating unit tests with descriptive names or: Would you name your children thing1 and thing2?,” in *Proceedings of the 26th ACM SIGSOFT International Symposium on Software Testing and Analysis*, pp. 57–67, ACM, 2017.

- [59] J. M. Rojas and G. Fraser, “Is search-based unit test generation research stuck in a local optimum?,” in *Proceedings of the 10th International Workshop on Search-Based Software Testing*, pp. 51–52, IEEE Press, 2017.
- [60] E. B. Swanson, “The dimensions of maintenance,” in *Proceedings of the 2nd international conference on Software engineering*, pp. 492–497, IEEE Computer Society Press, 1976.
- [61] S. C. B. de Souza, N. Anquetil, and K. M. de Oliveira, “A study of the documentation essential to software maintenance,” in *Proceedings of the 23rd annual international conference on Design of communication: documenting & designing for pervasive information*, pp. 68–75, ACM, 2005.
- [62] L. C. Briand, “Software documentation: how much is enough?,” in *Software Maintenance and Reengineering, 2003. Proceedings. Seventh European Conference on*, pp. 13–15, IEEE, 2003.
- [63] A. Forward and T. C. Lethbridge, “The relevance of software documentation, tools and technologies: a survey,” in *Proceedings of the 2002 ACM symposium on Document engineering*, pp. 26–33, ACM, 2002.
- [64] K. S. Jones, “Automatic summarising: The state of the art,” *Information Processing & Management*, vol. 43, no. 6, pp. 1449–1481, 2007.
- [65] N. Nazar, Y. Hu, and H. Jiang, “Summarizing software artifacts: A literature review,” *Journal of Computer Science and Technology*, vol. 31, no. 5, pp. 883–909, 2016.
- [66] G. Neubig, “Survey of methods to generate natural language from source code,” 2016.
- [67] X. Wang, Y. Peng, and B. Zhang, “Comment generation for source code: State of the art, challenges and opportunities,” *arXiv preprint arXiv:1802.02971*, 2018.
- [68] N. Dragan, M. L. Collard, and J. I. Maletic, “Reverse engineering method stereotypes,” in *Software Maintenance, 2006. ICSM’06. 22nd IEEE International Conference on*, pp. 24–34, IEEE, 2006.
- [69] N. Dragan, “Emergent laws of method and class stereotypes in object oriented software,” in *Software Maintenance (ICSM), 2011 27th IEEE International Conference on*, pp. 550–555, IEEE, 2011.
- [70] N. Dragan, M. L. Collard, and J. I. Maletic, “Using method stereotype distribution as a signature descriptor for software systems,” in *Software Maintenance, 2009. ICSM 2009. IEEE International Conference on*, pp. 567–570, IEEE, 2009.
- [71] L. Moreno and A. Marcus, “Jstereocode: automatically identifying method and class stereotypes in java code,” in *Proceedings of the 27th IEEE/ACM International Conference on Automated Software Engineering*, pp. 358–361, ACM, 2012.
- [72] R. P. Buse and W. R. Weimer, “Automatic documentation inference for exceptions,” in *Proceedings of the 2008 international symposium on Software testing and analysis*, pp. 273–282, ACM, 2008.

- [73] E. Hill, L. Pollock, and K. Vijay-Shanker, “Automatically capturing source code context of nl-queries for software maintenance and reuse,” in *Software Engineering, 2009. ICSE 2009. IEEE 31st International Conference on*, pp. 232–242, IEEE, 2009.
- [74] E. Hill, *Integrating natural language and program structure information to improve software search and exploration*. University of Delaware, 2010.
- [75] B. Liblit, A. Begel, and E. Sweetser, “Cognitive perspectives on the role of naming in computer programs,” in *Proceedings of the 18th annual psychology of programming workshop*, pp. 53–67, Citeseer, 2006.
- [76] G. Sridhara, E. Hill, D. Muppaneni, L. Pollock, and K. Vijay-Shanker, “Towards automatically generating summary comments for java methods,” in *Proceedings of the IEEE/ACM international conference on Automated software engineering*, pp. 43–52, ACM, 2010.
- [77] G. Sridhara, *Automatic generation of descriptive summary comments for methods in object-oriented programs*. University of Delaware, 2012.
- [78] P. W. McBurney and C. McMillan, “Automatic source code summarization of context for java methods,” *IEEE Transactions on Software Engineering*, vol. 42, no. 2, pp. 103–119, 2016.
- [79] G. Sridhara, L. Pollock, and K. Vijay-Shanker, “Automatically detecting and describing high level actions within methods,” in *Proceedings of the 33rd International Conference on Software Engineering*, pp. 101–110, ACM, 2011.
- [80] L. Moreno, J. Aponte, G. Sridhara, A. Marcus, L. Pollock, and K. Vijay-Shanker, “Automatic generation of natural language summaries for java classes,” in *Program Comprehension (ICPC), 2013 IEEE 21st International Conference on*, pp. 23–32, IEEE, 2013.
- [81] S. Rastkar, G. C. Murphy, and G. Murray, “Summarizing software artifacts: a case study of bug reports,” in *Proceedings of the 32nd ACM/IEEE International Conference on Software Engineering-Volume 1*, pp. 505–514, ACM, 2010.
- [82] N. Dragan, M. L. Collard, M. Hammad, and J. I. Maletic, “Using stereotypes to help characterize commits,” in *Software Maintenance (ICSM), 2011 27th IEEE International Conference on*, pp. 520–523, IEEE, 2011.
- [83] L. F. Cortés-Coy, M. Linares-Vásquez, J. Aponte, and D. Poshyvanyk, “On automatically generating commit messages via summarization of source code changes,” in *Source Code Analysis and Manipulation (SCAM), 2014 IEEE 14th International Working Conference on*, pp. 275–284, IEEE, 2014.
- [84] M. Linares-Vásquez, L. F. Cortés-Coy, J. Aponte, and D. Poshyvanyk, “Changscribe: A tool for automatically generating commit messages,” in *Software Engineering (ICSE), 2015 IEEE/ACM 37th IEEE International Conference on*, vol. 2, pp. 709–712, IEEE, 2015.
- [85] S. Jiang and C. McMillan, “Towards automatic generation of short summaries of commits,” in *Proceedings of the 25th International Conference on Program Comprehension*, pp. 320–323, IEEE Press, 2017.

- [86] S. Jiang, A. Armaly, and C. McMillan, “Automatically generating commit messages from diffs using neural machine translation,” in *Proceedings of the 32nd IEEE/ACM International Conference on Automated Software Engineering*, pp. 135–146, IEEE Press, 2017.
- [87] P. Loyola, E. Marrese-Taylor, and Y. Matsuo, “A neural architecture for generating natural language descriptions from source code changes,” *arXiv preprint arXiv:1704.04856*, 2017.
- [88] J. Shen, X. Sun, B. Li, H. Yang, and J. Hu, “On automatic summarization of what and why information in source code changes,” in *Computer Software and Applications Conference (COMPSAC), 2016 IEEE 40th Annual*, vol. 1, pp. 103–112, IEEE, 2016.
- [89] R. P. Buse and W. R. Weimer, “Automatically documenting program changes,” in *Proceedings of the IEEE/ACM international conference on Automated software engineering*, pp. 33–42, ACM, 2010.
- [90] J. Buckley, T. Mens, M. Zenger, A. Rashid, and G. Kniesel, “Towards a taxonomy of software change,” *Journal of Software: Evolution and Process*, vol. 17, no. 5, pp. 309–332, 2005.
- [91] S. Rastkar, G. C. Murphy, and A. W. Bradley, “Generating natural language summaries for crosscutting source code concerns,” in *Software Maintenance (ICSM), 2011 27th IEEE International Conference on*, pp. 103–112, IEEE, 2011.
- [92] M. Ghafari, C. Ghezzi, and K. Rubinov, “Automatically identifying focal methods under test in unit test cases,” in *Source Code Analysis and Manipulation (SCAM), 2015 IEEE 15th International Working Conference on*, pp. 61–70, IEEE, 2015.
- [93] B. Cornelissen, A. Van Deursen, L. Moonen, and A. Zaidman, “Visualizing testsuites to aid in software understanding,” in *Software Maintenance and Reengineering, 2007. CSMR’07. 11th European Conference on*, pp. 213–222, IEEE, 2007.
- [94] G. Meszaros, *xUnit test patterns: Refactoring test code*. Pearson Education, 2007.
- [95] M. Kamimura and G. C. Murphy, “Towards generating human-oriented summaries of unit test cases,” in *Program Comprehension (ICPC), 2013 IEEE 21st International Conference on*, pp. 215–218, IEEE, 2013.
- [96] K. Herbert, J. Goldhamer, and D. Ilvento, “Swummary: Self-documenting code,” 2016.
- [97] R. Jhala and R. Majumdar, “Path slicing,” in *ACM SIGPLAN Notices*, vol. 40, pp. 38–47, ACM, 2005.
- [98] B. Li, *Automatically Documenting Software Artifacts*. PhD thesis, College of William & Mary, 2018.
- [99] F. Beck, H. A. Siddiqui, A. Bergel, and D. Weiskopf, “Method execution reports: Generating text and visualization to describe program behavior,” in *Software Visualization (VISsOFT), 2017 IEEE Working Conference on*, pp. 1–10, IEEE, 2017.
- [100] M. E. Delamaro, J. C. Maldonado, and A. M. R. Vincenzi, “Proteum/im 2.0: An integrated mutation testing environment,” in *Mutation testing for the new century*, pp. 91–101, Springer, 2001.

- [101] B. S. Clegg, J. M. Rojas, and G. Fraser, “Teaching software testing concepts using a mutation testing game,” in *Software Engineering: Software Engineering Education and Training Track (ICSE-SEET)*, 2017 IEEE/ACM 39th International Conference on, pp. 33–36, IEEE, 2017.
- [102] A. Qusef, G. Bavota, R. Oliveto, A. De Lucia, and D. Binkley, “Scotch: Test-to-code traceability using slicing and conceptual coupling,” in *Software Maintenance (ICSM)*, 2011 27th IEEE International Conference on, pp. 63–72, IEEE, 2011.
- [103] B. Korel and J. Laski, “Dynamic program slicing,” *Information processing letters*, vol. 29, no. 3, pp. 155–163, 1988.
- [104] D. Poshyvanyk, A. Marcus, R. Ferenc, and T. Gyimóthy, “Using information retrieval based coupling measures for impact analysis,” *Empirical software engineering*, vol. 14, no. 1, pp. 5–32, 2009.
- [105] A. Leitner, M. Oriol, A. Zeller, I. Ciupa, and B. Meyer, “Efficient unit test case minimization,” in *Proceedings of the twenty-second IEEE/ACM international conference on Automated software engineering*, pp. 417–420, ACM, 2007.
- [106] A. Zeller, “Yesterday, my program worked. today, it does not. why?,” in *ACM SIGSOFT Software engineering notes*, vol. 24, pp. 253–267, Springer-Verlag, 1999.
- [107] M.-H. Liu, Y.-F. Gao, J.-H. Shan, J.-H. Liu, L. Zhang, and J.-S. Sun, “An approach to test data generation for killing multiple mutants,” in *Software Maintenance, 2006. ICSM’06. 22nd IEEE International Conference on*, pp. 113–122, IEEE, 2006.
- [108] J. Dolby, S. J. Fink, and M. Sridharan, “Tj watson libraries for analysis (wala),” *URL <http://wala.sf.net>*, 2015.
- [109] C. D. Newman, T. Sage, M. L. Collard, H. W. Alomari, and J. I. Maletic, “srcslice: a tool for efficient static forward slicing,” in *Software Engineering Companion (ICSE-C)*, IEEE/ACM International Conference on, pp. 621–624, IEEE, 2016.
- [110] A. Alali, H. Kagdi, and J. I. Maletic, “What’s a typical commit? a characterization of open source software repositories,” in *Program Comprehension, 2008. ICPC 2008. The 16th IEEE International Conference on*, pp. 182–191, IEEE, 2008.
- [111] L. P. Hattori and M. Lanza, “On the nature of commits,” in *Proceedings of the 23rd IEEE/ACM International Conference on Automated Software Engineering*, pp. III–63, IEEE Press, 2008.
- [112] S. Letovsky, “Cognitive processes in program comprehension,” *Journal of Systems and software*, vol. 7, no. 4, pp. 325–339, 1987.