

Optimising Continuous Integration using Test Case Prioritisation

Pieter De Clercq

Student number: 01503338

Supervisors: Prof. dr. Bruno Volckaert, Prof. dr. ir. Filip De Turck
Counsellors: Jasper Vaneessen, Dwight Kerkhove

Master's dissertation submitted in order to obtain the academic degree of
Master of Science in de informatica

Academic year 2019-2020

The author gives the permission to use this thesis for consultation and to copy parts of it for personal use. Every other use is subject to the copyright laws, more specifically the source must be extensively specified when using from this thesis.

De auteur geeft de toelating deze masterproef voor consultatie beschikbaar te stellen en delen van de masterproef te kopiëren voor persoonlijk gebruik. Elk ander gebruik valt onder de bepalingen van het auteursrecht, in het bijzonder met betrekking tot de verplichting de bron uitdrukkelijk te vermelden bij het aanhalen van resultaten uit deze masterproef.

Pieter De Clercq – May 21, 2020.

Acknowledgements

Completing this thesis would not have been possible without the help and support of many people, some of which I want to thank personally.

First of all, I want to thank prof. dr. Bruno Volckaert and prof. dr. ir. Filip De Turck for allowing me to propose this subject and for their prompt and clear responses to every question I have asked. I especially want to thank you for permitting me to insert a two-week hiatus during the Easter break, so I could help out on the UGent Dodona project.

Secondly, I want to express my gratitude towards my counsellors Jasper Vaneessen and Dwight Kerkhove, for steering me into researching this topic, as well as their guidance, availability, and willingness to review every intermediary version of this thesis.

Furthermore, I want to thank my parents, my brother Stijn and my family for convincing me and giving me the possibility to study at the university, to support me throughout my entire academic career and to provide me with the opportunity to pursue my childhood dreams.

Last, but surely not least, I want to thank my amazing friends, a few of them in particular. My best friend Robbe, for always being there when I need him even when I least expect it. For both supporting my wildest dreams while protecting me against my often unrealistic ideas and ambition to excel. Helena for never leaving my side, for always making me laugh when I don't want to, and most importantly to remind me that I should relax from time to time. Jana for my daily dose of laughter, fun and inexhaustible positivity. Tobiah for the endless design discussions and for outperforming me in almost every school project, to encourage me to continuously raise the bar and to never give up. Finally, I want to thank Doortje and Freija for answering my mathematical questions, regularly asking about my thesis progression and thereby motivating me to persevere.

Thank you.

Pieter – Ghent, 2020

Summary

Summary in English will come here.

Samenvatting

Nederlandse samenvatting komt hier.

Optimising Continuous Integration using Test Case Prioritisation

Pieter De Clercq

Supervisor(s): Prof. dr. B. Volckaert, Prof. dr. ir. F. De Turck, J. Vaneessen, D Kerkhove

Abstract—**This abstract is very abstract.**

Keywords—**words, will, appear, here, soon**

I. INTRODUCTIE

Things will appear here. [1]

REFERENCES

- [1] Michael Cusumano, Akindutire Michael, and Stanley Smith, "Beyond the waterfall : software development at microsoft," 02 1995.

Optimaliseren van Continue Integratie door middel van Test Prioritering

Pieter De Clercq

Supervisor(s): Prof. dr. B. Volckaert, Prof. dr. ir. F. De Turck, J. Vaneessen, D Kerkhove

Abstract—**Dit abstract is super abstract.**

Trefwoorden—**woorden, komen, hier**

I. INTRODUCTIE

Dingen komen hier. [1]

REFERENTIES

- [1] Michael Cusumano, Akindutire Michael, and Stanley Smith, "Beyond the waterfall : software development at microsoft," 02 1995.

Lay summary

Lay summary will come here.

Contents

Summary	iv
Summary (Dutch)	v
Extended abstract	vi
Extended abstract (Dutch)	vii
Lay summary	viii
1 Software Engineering [TODO REVISE]	2
1.1 Software Development Life Cycle	2
1.1.1 Test Suite Assessment	4
1.2 Agile Software Development	9
1.2.1 The need for Agile	9
1.2.2 Continuous Integration	9
2 Related work	13
2.1 Classification of approaches	15
2.1.1 Test Suite Minimisation	15
2.1.2 Test Case Selection	16
2.1.3 Test Case Prioritisation	16
2.2 Algorithms	18
2.2.1 Greedy algorithm	20
2.2.2 HGS	21
2.2.3 ROCKET algorithm	22
2.3 Adoption in testing frameworks	25
2.3.1 Gradle and JUnit	25
2.3.2 OpenClover	26
3 Proposed framework: VeloCity [TODO REVISE]	27
3.1 Design goals	27
3.2 Architecture	28
3.2.1 Agent	28
3.2.2 Controller	28
3.2.3 Predictor and Metrics	28
3.3 Pipeline	30
3.3.1 Initialisation	30

3.3.2	Prediction	31
3.3.3	Test case execution	33
3.3.4	Post-processing and analysis	33
3.4	Alpha algorithm	34
4	Evaluation	38
4.1	Test subjects	38
4.1.1	Dodona	38
4.1.2	Stratego	38
4.2	Research questions	39
4.3	Data collection	39
4.3.1	Travis CI build data	39
4.3.2	Dodona data	41
4.3.3	Stratego data	41
4.4	Results	42
4.4.1	RQ1: Probability of failure	42
4.4.2	RQ2: Average duration of a test run	42
4.4.3	RQ3: Consecutive failure probability	43
4.4.4	RQ4: Applying Test Case Prioritisation to Dodona	44
4.4.5	RQ5: Integrate VeloClty with Stratego	47
5	Conclusion	50
5.1	Future work	51
5.1.1	Java agent	51
5.1.2	Predictions	51
5.1.3	Meta predictor	52
5.1.4	Final enhancements	53
	Appendices	59

Glossary

CI Continuous Integration. 13

MapReduce a programming paradigm that allows large amounts of data to be processed in a distributed manner. 40

TCP Test Case Prioritisation. 14

TCS Test Case Selection. 14

TSM Test Suite Minimisation. 14, 15

Chapter 1

Software Engineering [TODO REVISE]

The Institute of Electrical and Electronics Engineers [IEEE] defines the practice of Software Engineering as: "Application of a systematic, disciplined, quantifiable approach to the development, operation and maintenance of software; that is, the application of engineering to software" [19, p. 421]. The word "systematic" in this definition, emphasises the need for a structured process, depicting guidelines and models that describe how software should be developed the most efficient way possible. Such a process does exist and it is often referred to as the Software Development Life Cycle (SDLC) [19, p. 420]. In the absence of a model, i.e. when the developer does what they deem correct without following any rules, the term *Cowboy coding* is used [21, p. 34].

1.1 Software Development Life Cycle

An implementation of the SDLC consists of two major components. First, the process is broken down into several smaller phases. Depending on the nature of the software, it is possible to omit steps or add more steps. I have compiled a simple yet generic approach from multiple sources [13, 18], to which most software projects adhere. This approach consists of five phases.

1. **Requirements phase:** This is the initial phase of the development process. During this phase, the developer gets acquainted with the project and compiles a list of the desired functionalities [18]. Using this information, the developer eventually decides on the required hardware specifications and possible external software which will need to be acquired.
2. **Design phase:** After the developer has gained sufficient knowledge about the project requirements, they can use this information to draw an architectural design of the application. This design consists of multiple documents, including user stories and UML-diagrams.
3. **Implementation phase:** During this phase, the developer will write code according to the specifications defined in the architectural designs.
4. **Testing phase:** This is the most important phase. During this phase, the implementation is tested to identify potential bugs before the application is used by other users.

5. **Operational phase:** In the final phase, the project is fully completed and it is integrated in the existing business environment.

Subsequently, a model is chosen to define how to transition from one phase into another phase. A manifold of models exist [13], each having advantages and disadvantages, but I will consider the basic yet most widely used model, which is the Waterfall model by Benington [4]. The initial Waterfall model required every phase to be executed sequentially and in order, cascading. However, this imposes several issues, the most prevalent being the inability to revise design decisions taken in the second phase, when performing the actual implementation in the third phase. To mitigate this, an improved version of the Waterfall model was proposed by Royce [31]. This version allows a phase to transition back to a previous phase (Figure 1.1).

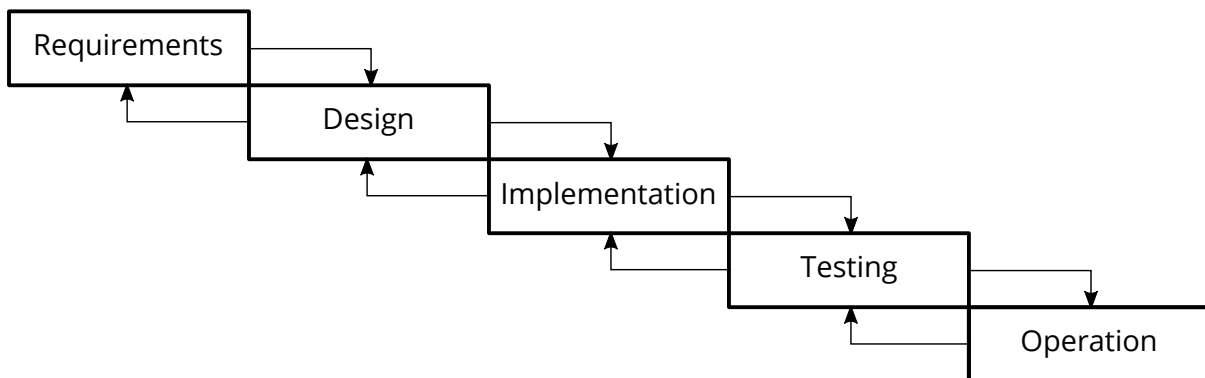


Figure 1.1: Improved Waterfall model by Royce

In this thesis I will solely focus on the implementation and testing phase, as these are the most time-consuming phases of the entire process. The modification to the Waterfall model by Royce is particularly useful when applied to these two phases, in the context of *software regressions*. A regression [28] is a feature that was previously working correctly, but is now malfunctioning. This behaviour can have external causes, such as a change in the system clock because of daylight saving time, but can also be the result of a change to another, seemingly unrelated part of the application code [17].

Software regressions and other functional bugs can ultimately incur disastrous effects, such as severe financial loss or damage to the reputation of the software company. The most famous example in history is without any doubt the explosion of the Ariane 5-rocket, which was caused by an integer overflow [22]. In order to reduce the risk of bugs, malfunctioning components should be detected as soon as possible to proactively defend against potential failures. Because of this reason, the testing phase is to be considered as the most important phase of the entire development process and an application should therefore include sufficient tests. The collection of all tests included in an application, or a smaller chosen subset of certain tests, is referred to

as the *test suite*. Tests can be classified in multiple categories, this thesis will consider three distinguishable categories:

1. **Unit test:** This is the most basic kind of test. The purpose of a unit test is to verify the behaviour of an individual component [35]. The scope of a unit test should be limited to a small and isolated piece of code, such as one function. Unit tests are typically implemented as *white-box tests* [17, p. 12]. A white-box test is constructed by manually inspecting the function under test, to identify important *edge values*. The unit test should then feed these values as arguments to the function under test, to observe its behaviour. Common edge cases include zero, negative numbers, empty arrays or array boundaries that might result in an overflow.
2. **Integration test:** A more advanced test, an integration test verifies the interaction between multiple individually tested components [35]. Examples of integration tests include the communication between the front-end and the back-end side of an application. As opposed to unit tests, an integration test is an example of a *black-box test* [17, p. 6], meaning that implementation-specific details should be irrelevant or unknown when writing an integration test.
3. **Regression test:** After a regression has been detected, a regression test [19, p. 372] is added to the test suite. This regression test should replicate the exact conditions and sequence of actions that have caused the regression, to warn the implementation against subsequent failures if the same conditions would reapply in the future.

1.1.1 Test Suite Assessment

Coverage

The most frequently used metric to measure the quantity and thoroughness of a test suite is the *code coverage* or *test coverage* [19, p. 467]. The test coverage is expressed as a percentage and indicates which fraction of the application code is affected by code in the test suite. Internally, this works by augmenting every statement in the application code using binary instrumentation. A hook is inserted before and after every statement to keep track of which statements are executed during tests. Many different criteria exist to interpret these instrumentation results and thus to express the fraction of covered code [27], the most commonly used ones are *statement coverage* and *branch coverage*.

Statement coverage expresses the fraction of code statements that are executed in any test of the test suite [17], out of all executable statements in the application code. Analogously, the fraction of lines covered by a test may be used to calculate the *line coverage* percentage. Since one statement can span multiple lines and one line may also contain more than one statement, both of these criteria implicitly represent the same value. Statement coverage is heavily criticised in literature [27, p. 37], since it is possible to achieve a statement coverage percentage of 100% on a code fragment which can be proven to be incorrect. Consider the code fragment in Listing 1.1. If a test would call the `example`-function with arguments $\{a = 1, b = 2\}$, the test will pass and every statement will be covered, resulting in a statement coverage of 100%. However, it is clear to see that if the function would be called with arguments $\{a = 0, b = 0\}$, a *division-by-zero* error would be raised, resulting in a crash. This very short example already indicates that statement coverage is not trustworthy, yet it may still be useful for other purposes, such as detecting unreachable code which may safely be removed.

```
1 int example(int a, int b) {  
2     if (a == 0 || b != 0) {  
3         return a / b;  
4     }  
5 }
```

Listing 1.1: Example of irrelevant statement coverage in C.

Branch coverage on the other hand, requires that every branch of a conditional statement is traversed at least once [27, p. 37]. For an `if`-statement, this results in two tests being required, one for every possible outcome of the condition (`true` or `false`). For a `loop`-statement, this requires a test case in which the loop body is never executed and another test case in which the loop body is always executed. Remark that while this criterion is stronger than statement coverage, it is still not sufficiently strong to detect the bug in Listing 1.1. In order to mitigate this, *multiple-condition coverage* [27, p. 40] is used. This criterion requires that for every conditional statement, every possible combination of subexpressions is evaluated at least once. Applied to Listing 1.1, the `if`-statement is only covered if the following four cases are tested, which is sufficient to detect the bug.

- $a = 0, b = 0$
- $a = 0, b \neq 0$
- $a \neq 0, b = 0$
- $a \neq 0, b \neq 0$

It should be self-evident that achieving and maintaining a coverage percentage of 100% at all times is critical. However, this does not necessarily imply that all lines, statements or branches need to be covered explicitly [7]. Some parts of the code might simply be irrelevant or untestable. Examples include wrapper or delegation methods that simply call a library function. All major programming languages have frameworks and libraries available to collect coverage information during test execution, and each of these frameworks allows the developer to exclude parts of the code from the final coverage calculation. As of today, the most popular options are JaCoCo¹ for Java, coverage.py² for Python and simplecov³ for Ruby. These frameworks are able to generate in-depth statistics on which parts of the code are covered and which parts require more tests, as illustrated in Figure 1.3.

Mutation testing

Whereas code coverage can be used to identify whether or not a part of the code is currently affected by the test suite, *mutation testing* can be used to measure its quality and ability to detect future failures. This technique creates several syntactically different instances of the source code, referred to as *mutants*. A mutant can be created by applying one or more *mutation operators* to the original source code. These mutation operators are aimed at simulating typical mistakes that developers tend to make, such as the introduction of off-by-one errors, removal of statements and replacement of logical connectors [30]. The *mutation order* refers to the amount of mutation operators that have been applied consecutively to an instance of the code. This order is traditionally rather low, as a result of the *Competent Programmer Hypothesis*, which states that programmers develop programs which are near-correct [20].

Creating and evaluating the mutant versions of the code is a computationally expensive process and requires human intervention, which is why very few software developers have managed to employ this technique in practice. Figure 1.2 shows how mutation testing is performed. First of all, the mutation system takes the original program P and a set of test cases T . Then, several mutation operators are applied to construct a large set of mutants P' . The next step is to evaluate every test case t on the original program P to verify its correctness, this is a task that needs to be performed manually. If at least one of these test cases proves incorrect, a bug has been found in the original program, which needs to be resolved before the mutation analysis can continue. When P successfully passes every test case, every test case are evaluated for each of the mutants. A mutant p' is said to be “killed” if its output is different from

¹<https://www.jacoco.org/jacoco/>

²<https://github.com/nedbat/coveragepy>

³<https://github.com/colszowka/simplecov>

P for at least one test case, otherwise it is considered “surviving”. After executing all test cases, the set of surviving mutants should be analysed in order to introduce subsequent test cases that can be used to kill them. However, it is also possible that the surviving mutants are functionally equivalent to P . This needs to be verified manually, since the detection of program equivalence is impossible [20, 30].

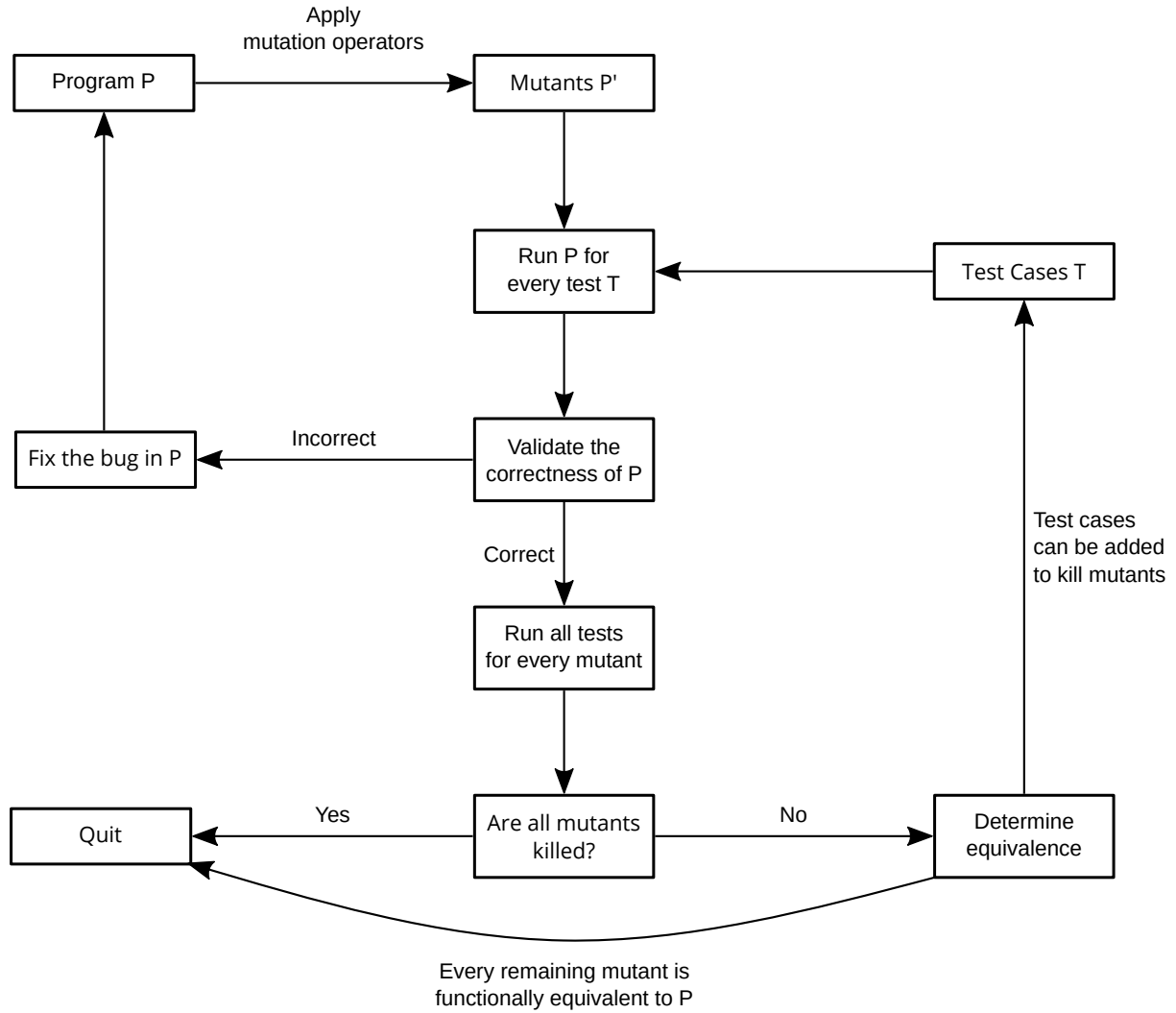





Figure 1.2: Process of Mutation Testing (based on [30])

After every mutant has either been killed or marked equivalent to the original problem, the test suite is assigned a *mutation score* which is calculated using Equation 1.1. In an ideal test suite, this score should be equal to 1, indicating that the test suite was able to detect every mutant.

$$\text{Mutant Score} = \frac{\text{killed mutants}}{\text{non-equivalent mutants}} \quad (1.1)$$

io.github.thepieterdc.http.impl

Element	Missed Instructions	Cov.	Missed Branches	Cov.	Missed	Cxty	Missed	Lines	Missed	Methods	Missed	Classes
HttpClientImpl		59%		14%	7	14	18	40	2	9	0	1
HttpResponseImpl		55%	n/a	n/a	9	15	10	22	9	15	0	1
Total	88 of 211	58%	6 of 7	14%	16	29	28	62	11	24	0	2

(a) JaCoCo coverage report of <https://github.com/thepieterdc/dodona-api-java>**Coverage report: 75%**

Module ↓	statements	missing	excluded	coverage
awesome/__init__.py	4	1	0	75%
<pre> 1 def smile(): 2 return ":" 3 4 def frown(): 5 return ":(</pre>				
Total	4	1	0	75%

(b) coverage.py report of <https://github.com/codecov/example-python>**Helpers (88.41% covered at 22.84 hits/line)**

12 files in total. 716 relevant lines. 633 lines covered and 83 lines missed

File	% covered	Lines	Relevant Lines	Lines covered	Lines missed	Avg. Hits / Line
app/helpers/standard_form_builder.rb	100.0 %	5	3	3	0	11.0
app/helpers/renderers/feedback_code_renderer.rb	100.0 %	25	16	16	0	5.4
app/helpers/institutions_helper.rb	100.0 %	2	1	1	0	1.0
app/helpers/api_tokens_controller_helper.rb	100.0 %	2	1	1	0	1.0
app/helpers/renderers/pythia_renderer.rb	93.94 %	290	165	155	10	3.6
app/helpers/renderers/feedback_table_renderer.rb	90.59 %	349	202	183	19	16.8
app/helpers/exercise_helper.rb	90.16 %	125	61	55	6	3.5
app/helpers/courses_helper.rb	86.67 %	36	15	13	2	28.4
app/helpers/repository_helper.rb	85.71 %	11	7	6	1	2.6
app/helpers/application_helper.rb	85.59 %	220	111	95	16	62.6
app/helpers/users_helper.rb	84.62 %	20	13	11	2	1.4
app/helpers/renderers/lcs_html_differ.rb	77.69 %	236	121	94	27	38.2

Showing 1 to 12 of 12 entries

(c) simplecov report of <https://github.com/dodona-edu/dodona>

Figure 1.3: Statistics from Code coverage tools

1.2 Agile Software Development

1.2.1 The need for Agile

In the wake of the world economic crisis, software companies were forced to devote efforts into researching how their overall expenses could be reduced. This research has concluded that in order to reduce financial risks, the *time-to-market* of an application should be as short as possible. In order to accomplish this, further research was conducted, resulting in an increase of attention for agile methodologies in scientific literature [16]. As was previously described in ??, agile methodologies strive to deliver a minimal version as soon as possible, allowing additional functionality to be added in an incremental fashion. This effectively results in a shorter *time-to-market* and lower costs, since the company can decide to cancel the project much earlier in the process.

In addition to a reduced time-to-market, maintaining an agile workflow has also proven beneficial to the success rate of development. A study performed by The Standish Group revealed that the success rate of agile projects is more than three times higher compared to when traditional methodologies are practised, as illustrated in Figure 1.4.

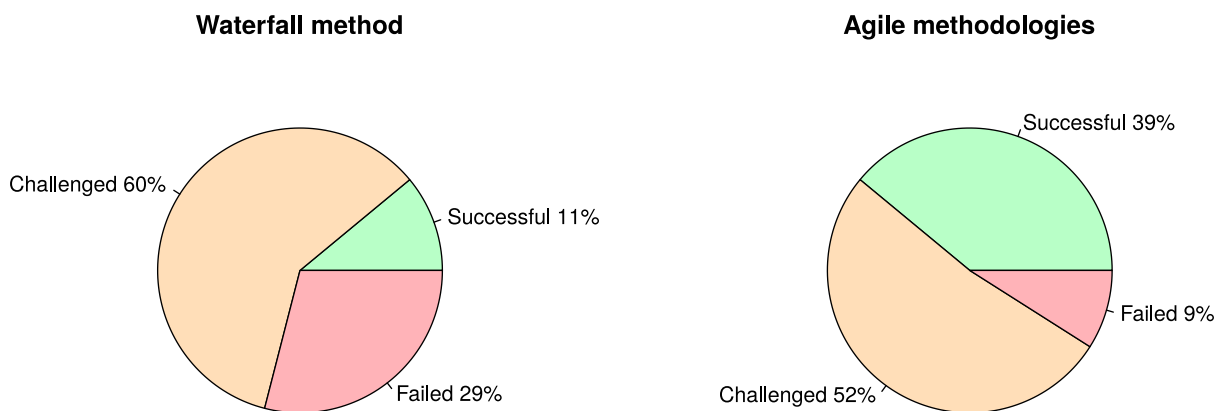


Figure 1.4: Success rate of Agile methodologies [14].

1.2.2 Continuous Integration

In traditional software development, the design phase typically leads to a representation of the required functionality in multiple, stand-alone modules. Subsequently, every module is implemented separately by individual developers. Afterwards, an attempt is made to integrate all the modules into the final application, an event to which Meyer refers to as the “Big Bang” [26, p.103]. The name *Big Bang* reflects the complex nature of this operation. This can prove to be a challenging operation, because

every developer can take unexpected assumptions at the start of the project, which may ultimately result in mutually incompatible components. Furthermore, since the code was written over a span of several weeks to months, the developers often need to rewrite code that they have not touched in a long time. Eventually this will lead to unanticipated delays and costs [32].

Contrarily, agile development methodologies advocate the idea of frequent, yet small deliveries (??). Consequently, this implies that the code is built often and that the modules are integrated multiple times, on a *continuous* basis, rather than just once at the end, thus allowing for early identification of problems [12]. This practice of frequent builds is referred to as *Continuous Integration* [25, 26]. It should be noted that this idea has existed and has been applied before the agile manifesto was written. The first notorious software company that has adopted this practice is Microsoft, already in 1989 [6, p.11]. Cusumano reports that Microsoft typically builds the entire application at least once per day [6, p.12], therefore requiring developers to integrate and test their changes multiple times per day.

The introduction of Continuous Integration [CI] in software development has important consequences on the life cycle. Where the waterfall model used a cascading life cycle, Continuous Integration employs a circular, repetitive structure consisting of three phases, as visualised in Figure 1.5.

1. **Implementation:** In the first phase, every developer individually writes code for the module they were assigned to. At a regular interval, the code is committed to the remote repository.
2. **Integration:** When the code is committed, the developer simultaneously fetches the changes to other modules. Afterwards, the developer must integrate the changes with his own module, to ensure compatibility. In case a conflict occurs, the developer is responsible for its resolution [25].
3. **Test:** After the module has successfully been integrated, the test suite is run to ensure no bugs have been introduced.

Adopting Continuous Integration can prove to be a lengthy and repetitive task. Luckily, a variety of tools and frameworks exist to automate this process. Essentially, these tools are typically attached to a version control system (e.g. Git, Mercurial, ...), using a `post-receive` hook. Every time a commit is pushed by one of the developers, the CI system is notified, after which the code is automatically built and tests are executed. Optionally, the system can be configured to automatically publish successful

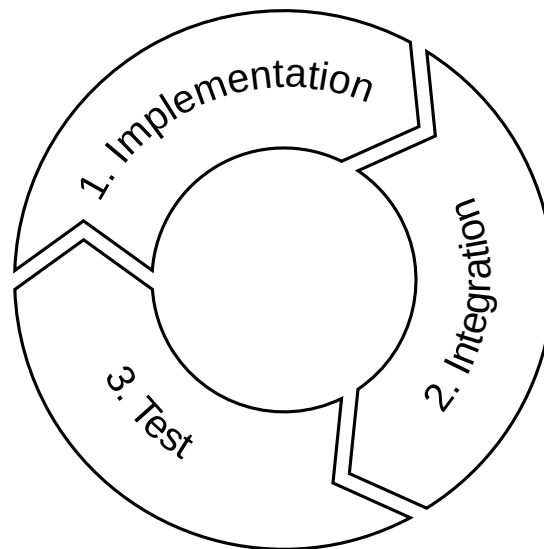


Figure 1.5: Development Life Cycle with Continuous Integration

runs to the end users, a process referred to as *Continuous Delivery*. I will now proceed by discussing four prominent Continuous Integration systems.

Jenkins

Jenkins CI⁴ was started as a hobby project in 2004 by Kohsuke Kawaguchi, a former employee of Sun Microsystems. Jenkins is programmed in Java and profiles itself as “The leading open source automation server”. It was initially named Hudson, but after Sun was acquired by Oracle, issues related to the trademark Hudson arose. In response, the developer community decided to migrate the Hudson code to a new repository and rename the project to Jenkins [32]. As of today, Jenkins is still widely used for many reasons. Since it is open source and its source code is located on GitHub, it is free to use and can be self-hosted by the developers in a private environment. Furthermore, Jenkins provides an open ecosystem to support developers into writing new plugins and extending its functionality. A market research conducted by ZeroTurnaround in 2016 revealed that Jenkins is the preferred Continuous Integration tool by 60% of the developers [23].



Figure 1.6: Logo of Jenkins CI (<https://jenkins.io/>)

⁴<https://jenkins.io/>

GitHub Actions

Following the successful beta of GitHub Actions which had started in August 2019, GitHub launched its own Continuous Integration system later that year in November⁵. GitHub Actions executes builds in the cloud on servers owned by GitHub and can therefore only be used in conjunction with a GitHub repository, support for GitHub Enterprise repositories is not yet available. The developers can define builds using *workflows* that can be configured to run both on Linux, Windows as well as OSX hosts. Private repositories are allowed a fixed amount of free build minutes per month, while builds of public repositories are always free of charges [9]. Similar to Jenkins, GitHub Actions can be extended with custom plugins. These plugins can be created using either a Docker container, or in native JavaScript, which allows faster execution [1]. It should be noted however that due to the recent nature of this system, not many plugins have been created yet.



Figure 1.7: Logo of GitHub Actions (<https://github.com/features/actions>)

GitLab CI

GitLab, the main competitor of GitHub, announced its own Continuous Integration service in late 2012 named GitLab CI⁶. GitLab CI builds are configured using *pipelines* and are executed by *GitLab Runners*. These runners are operated by developers on their own infrastructure. Additionally, GitLab also offers the possibility to use *shared runners*, which are hosted by themselves [2]. Equivalent to the aforementioned GitHub Actions, shared runners can be used for free by public repositories and are bounded by quota for private repositories [11]. A downside of using GitLab CI is the lack of a community-driven plugin system, however this is a planned feature⁷.



Figure 1.8: Logo of GitLab CI (<https://gitlab.com/>)

⁵<https://github.blog/2019-08-08-github-actions-now-supports-ci-cd/>

⁶<https://about.gitlab.com/blog/2012/11/13/continuous-integration-server-from-gitlab/>

⁷<https://gitlab.com/gitlab-org/gitlab/issues/15067>

Travis CI

The final Continuous Integration platform which I will discuss is Travis CI. This Continuous Integration system was launched in 2011 and can only be used in addition to an existing GitHub repository. Travis CI build tasks can be configured in a similar fashion as GitLab CI, but the builds can exclusively be executed on their servers. Besides running builds after a commit has been pushed to the repository, it is also possible to schedule daily, weekly or monthly builds using cron jobs. Similar to GitHub Actions, open-source projects can be built at zero cost and a paid plan exists for private repositories [8]. It is not possible to create custom plugins, however Travis CI already features built-in support for a variety of programming languages. In 2020, almost 1 million projects are being built using Travis CI [33].



Figure 1.9: Logo of Travis CI (<https://travis-ci.com/>)

Chapter 2

Related work

In the previous chapter, we have stressed the paramount importance of frequently integrating one's changes into the upstream repository. This process can prove to be a complex and lengthy operation. As a result, software engineers have sought and found ways to automate this task. These solutions and practices embody Continuous Integration (CI). However, CI is not the golden bullet for software engineering, as there is a flip side to applying this practice. After every integration, we must execute the entire test suite to ensure that we have not introduced any regressions. As the project evolves and the size of the codebase increases, the number of test cases will increase accordingly to preserve a sufficiently high coverage level [29]. Walcott, Soffa and Kapfhammer illustrate the magnitude of this problem by providing an example of a project consisting of 20 000 lines of code, whose test suite requires up to seven weeks to complete [34].

Fortunately, developers and researchers have found multiple techniques to address

the scalability issues of ever-growing test suites. We can classify the techniques currently known in literature into three categories [29]. These categories are Test Suite Minimisation (TSM), Test Case Selection (TCS) or Test Case Prioritisation (TCP). We can apply each technique to every test suite, but the outcome will be different. TSM and TCS will have an impact on the execution time of the test suite, at the cost of a reduced test coverage level. In contrast, TCP will have a weaker impact on the execution time but will not affect the test adequacy.

The following sections will discuss these three approaches in more detail and provide accompanying algorithms. Because the techniques are very similar, the corresponding algorithms can (albeit with minor modifications) be used interchangeably for every approach. The final section of this chapter will investigate the adoption and integration of these techniques in modern software testing frameworks.

2.1 Classification of approaches

2.1.1 Test Suite Minimisation

The first technique is called Test Suite Minimisation, also referred to as *Test Suite Reduction* in literature. This technique will try to reduce the size of the test suite by permanently removing redundant test cases. This problem has been formally defined by Rothermel in definition 1 [36] and illustrated in Figure 2.1.

ALLES HIERONDER NOG DOEN

Definition 1 (Test Suite Minimisation).

Given:

- $T = \{t_1, \dots, t_n\}$ a test suite consisting of test cases t_i .
- $R = \{r_1, \dots, r_n\}$ a set of requirements that must be satisfied in order to provide the desired “adequate” testing of the program.
- $G = \{T_1, \dots, T_n\} \subseteq T$ subsets of test cases, one associated with each of the requirements r_i , such that any one of the test cases $t_j \in T_i$ can be used to satisfy requirement r_i .

Test Suite Minimisation is then defined as the task of finding a subsetset T' of test cases $t_j \in T$ that satisfies all requirements r_i .

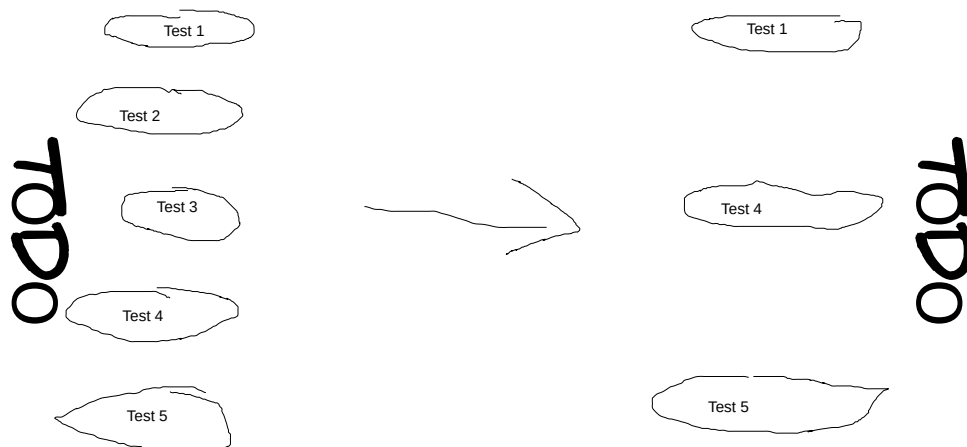


Figure 2.1: Test Suite Minimisation

If we apply this definition to the concepts introduced in ??, the requirements R can be interpreted as lines in the codebase that must be covered. With respect to the definition, a requirement can be satisfied by any test t_j that belongs to subset T_i of T . Observe that the problem of finding T' is closely related to the *hitting set problem* (2) [36].

Definition 2 (Hitting Set Problem).*Given:*

- $S = \{s_1, \dots, s_n\}$ a finite set of elements.
- $C = \{c_1, \dots, c_n\}$ a collection of sets, with $\forall c_i \in C : c_i \subseteq S$.
- K a positive integer, $K \leq |S|$.

The hitting set is a subset $S' \subseteq S$ such that S' contains at least one element from each subset in C .

In the context of Test Suite Minimisation, T' is precisely the hitting set of T_i s. In order to effectively minimise the amount of tests in the test suite, T' should be the minimal hitting set [36], which is an NP-complete problem as it can be reduced to the *Vertex Cover*-problem [10].

2.1.2 Test Case Selection

The second algorithm closely resembles the previous one. Instead of determining the minimal hitting set of the test suite in order to permanently remove tests, this algorithm has a notion of context. Prior to the execution of the tests, the algorithm performs a *white-box static analysis* of the codebase to identify which parts have been changed. Subsequently, only the tests regarding modified parts are executed, making the selection temporary (Figure 2.2) and modification-aware [36]. Rothermel and Harrold define this formally in 3.

Definition 3 (Test Case Selection).*Given:*

- P the previous version of the codebase
- P' the current (modified) version of the codebase
- T the test suite

Test Case Selection aims to find a subset $T' \subseteq T$ that is used to test P' .

2.1.3 Test Case Prioritisation

Where the previous algorithms both attempted to execute as few tests as possible, it might sometimes be desired or even required that all tests pass. In this case, the previous ideas can be used as well. In Test Case Prioritisation, we want to find a permutation of the sequence of all tests instead of eliminating certain tests. The order

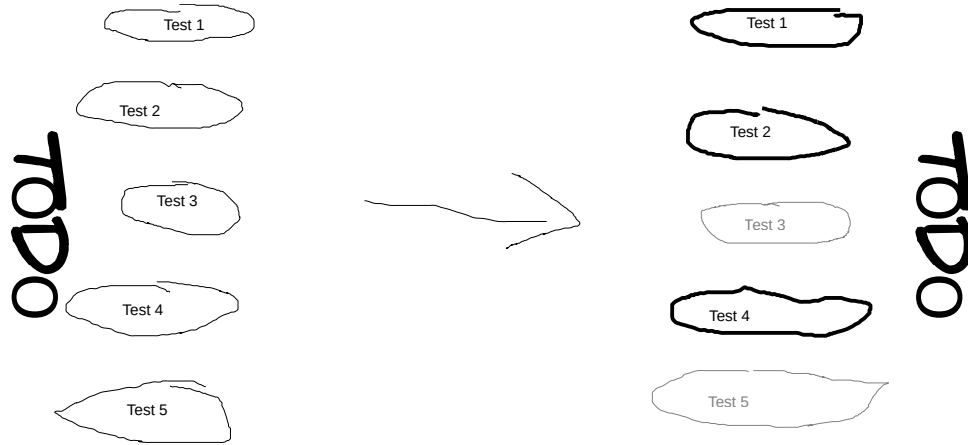


Figure 2.2: Test Case Selection

of the permutation is chosen specifically to achieve a given goal as soon as possible, allowing for early termination of the test suite upon failure [36]. Some examples of goals include covering as many lines of code as fast as possible, or early execution of tests with a high probability of failure. A formal definition of this algorithm is provided in 4.

Definition 4 (Test Case Prioritisation).

Given:

- T the test suite
- PT the set of permutations of T
- $f : PT \mapsto \mathbb{R}$ a function from a subset to a real number, this function is used to compare sequences of tests to find the optimal permutation.

Test Case Prioritisation finds a permutation $T' \in PT$ such that $\forall T'' \in PT : f(T') \geq f(T'') \Rightarrow (T'' \neq T')$

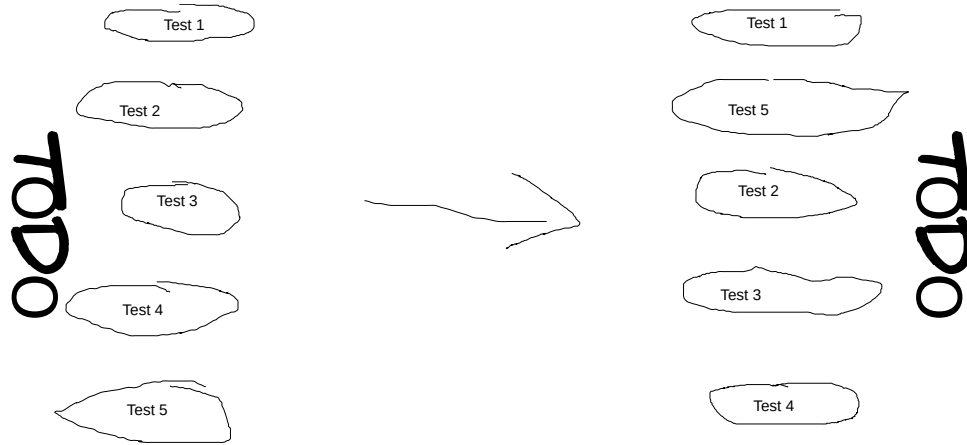


Figure 2.3: Test Case Prioritisation

2.2 Algorithms

In subsection 2.1.1 the relation was explained between applying Test Suite Minimisation and finding the minimal hitting set of the test suite and the set of requirements, which is an NP-complete problem. Therefore, the use of *heuristics* is required. A heuristic is an experience-based method that can be used to solve a hard to compute problem by finding a fast approximation [17]. However, the found solution will mostly be suboptimal or might sometimes even fail to find any solution at all. Considering its relation to the minimal hitting set problem, heuristics that are known to literature for solving this problem can also be used to implement Test Suite Minimisation. A selection of these heuristics will be discussed below. It should be noted however that the used terminology and naming of the variables might have been changed to ensure mutual consistency. Every algorithm has been adapted to adhere to the conventions provided in 5 and 6.

Definition 5 (Naming convention).

- C : the set of all lines in the application source code that are covered by at least one test case $t \in TS$.
 - CT_l denotes the test group l , which corresponds to the set of all tests $t \in TS$ that cover source code line $l \in C$.
- RS : the representative set of test cases, these are the test cases that have been selected by the algorithm.
- TS : the set of all test cases in the test suite.
 - TL_t denotes the set of all source code lines that are covered by test $t \in TS$.

Definition 6 (Cardinality). *For a finite set S , the cardinality $|S|$ is defined as the number of elements in S . In case of potential confusion, the cardinality of S can also be denoted as $\text{Card}(S)$.*

2.2.1 Greedy algorithm

The first algorithm is a *greedy* heuristic, which was originally designed by Chvatal to find an approximation for the set-covering problem [29]. A greedy algorithm always makes a locally optimal choice, assuming that this will eventually lead to a globally optimal solution [5]. Algorithm 1 presents the Greedy algorithm for Test Suite Minimisation. The goal of the algorithm is to construct a set of test cases that cover every line in the code, by requiring as few tests as possible.

Initially, the algorithm starts with an empty result set RS , the set TS of all test cases and the set C of all coverable source code lines. Furthermore, TL_t denotes the set of source code lines in C that are covered by test case $t \in TS$. Subsequently, the algorithm iteratively selects test cases from TS and adds them to RS . The locally optimal choice is to always select the test case that will contribute the most still uncovered lines, ergo the test t for which the cardinality of the intersection between C and TL_t is maximal. After every iteration, the now covered lines TL_t are removed from C and the selection process is repeated until C is empty. Upon running the tests, only the tests in RS must be executed. This algorithm can be converted to make it applicable to Test Case Prioritisation by converting the set RS to a list to maintain the order in which the test cases were selected, which is equivalent to the prioritised order of execution.

Algorithm 1 Greedy algorithm for Test Suite Minimisation

```

1: Input: Set  $TS$  of all test cases,
      Set  $C$  of all source code lines that are covered by any  $t \in TS$ ,
       $TL_t$  the set of all lines are covered by test case  $t \in TS$ .
2: Output: Subset  $RS \subseteq TS$  of tests to execute.
3:  $RS \leftarrow \emptyset$ 
4: while  $C \neq \emptyset$  do
5:    $t_{max} \leftarrow 0$ 
6:    $tl_{max} \leftarrow \emptyset$ 
7:   for all  $t \in TS$  do
8:      $tl_{current} \leftarrow C \cap TL_t$ 
9:     if  $|tl_{current}| > |tl_{max}|$  then
10:       $t_{max} \leftarrow t$ 
11:       $tl_{max} \leftarrow tl_{current}$ 
12:    $RS \leftarrow RS \cup \{t_{max}\}$ 
13:    $C \leftarrow C \setminus tl_{max}$ 

```

2.2.2 HGS

The second algorithm was created by Harrold, Gupta and Soffa [15]. This algorithm constructs the minimal hitting set of the test suite in an iterative fashion. As opposed to the greedy algorithm (subsection 2.2.1), the HGS algorithm considers the test groups CT instead of the set TLt to obtain a list of test cases that cover all source code lines. More specifically, this algorithm considers the distinct test groups, denoted as CTD . Two test groups are considered indistinct if they differ in at least one test case. The pseudocode for this algorithm is provided in Algorithm 2.

Similar to the previous algorithm, an empty representative set RS is constructed in which the selected test cases will be stored. The algorithm begins by iterating over every source code line $l \in C$ and constructing the corresponding set of test groups CT_l . As mentioned before, for performance reasons this set is reduced to CTD , only retaining distinct test groups. Next, the algorithm selects every test group of which the cardinality is equal to 1 and adds these to RS . This corresponds to every test case that covers a line of code, which is exclusively covered by that single test case. Subsequently, the lines that are covered by any of the selected test cases are removed from C . This process is repeated for an incremented cardinality, until every line in C is covered. Since the remaining test groups will now contain more than one test case, the algorithm needs to make a choice on which test case to select. The authors have chosen that the test case that occurs in the most test groups is preferred. In the event of a tie, this choice is deferred until the next iteration.

The authors have provided an accompanying calculation of the computational time complexity of this algorithm [15]. With respect to the naming convention introduced in 5, additionally let n denote the number of distinct test groups CTD , nt the number of test cases $t \in TS$ and MAX_CARD the cardinality of the largest test group. The HGS algorithm consists of two steps which are performed repeatedly. The first step involves computing the number of occurrences of every test case t in each test group. Given that there are n distinct test groups and, in the worst case scenario, each test group can contain MAX_CARD test cases which all need to be examined once, the computational cost of this step is equal to $O(n * MAX_CARD)$. In order to determine which test case should be included in the representative set RS , the algorithm needs to find all test cases for which the number of occurrences in all test groups is maximal, which requires at most $O(nt * MAX_CARD)$. Since every repetition of these two steps adds a test case that belongs to at least one out of n test groups to the representative set, the overall runtime of the algorithm is $O(n * (n + nt) * MAX_CARD)$.

Algorithm 2 HGS algorithm ([15])

```

1: Input: Distinct test groups  $T_1, \dots, T_n \in CDT$ , containing test cases from  $TS$ .
2: Output: Subset  $RS \subseteq TS$  of tests to execute.
3:  $marked \leftarrow \text{array}[1 \dots n]$  ▷ initially false
4:  $MAX\_CARD \leftarrow \max\{Card(T_i) | T_i \in CDT\}$ 
5:  $RS \leftarrow \bigcup \{T_i | Card(T_i) = 1\}$ 
6: for all  $T_i \in CDT$  do
7:   if  $T_i \cap RS \neq \emptyset$  then  $marked[i] \leftarrow true$ 
8:  $current \leftarrow 1$ 
9: while  $current < MAX\_CARD$  do
10:   $current \leftarrow current + 1$ 
11:  while  $\exists T_i : Card(T_i) = current, marked[i] = false$  do
12:     $list \leftarrow \{t | t \in T_i : Card(T_i) = current, marked[i] = false\}$ 
13:     $next \leftarrow SelectTest(current, list)$ 
14:     $reduce \leftarrow false$ 
15:    for all  $T_i \in CDT$  do
16:      if  $next \in T_i$  then
17:         $marked[i] = true$ 
18:        if  $Card(T_1) = MAX\_CARD$  then  $reduce \leftarrow true$ 
19:      if  $reduce$  then
20:         $MAX\_CARD \leftarrow \max\{Card(T_i) | marked[i] = false\}$ 
21:       $RS \leftarrow RS \cup \{next\}$ 
22: function  $SELECTTEST(size, list)$ 
23:   $count \leftarrow \text{array}[1 \dots nt]$ 
24:  for all  $t \in list$  do
25:     $count[t] \leftarrow |\{T_j | t \in T_j, marked[T_j] = false, Card(T_j) = size\}|$ 
26:   $tests \leftarrow \{t | t \in list, count[t] = \max(count)\}$ 
27:  if  $|tests| = 1$  then return  $tests[0]$ 
28:  else if  $|tests| = MAX\_CARD$  then return  $tests[0]$ 
29:  else return  $SelectTest(size + 1, tests)$ 

```

2.2.3 ROCKET algorithm

In contrast to the previously discussed algorithms which focused on Test Suite Minimisation, the ROCKET algorithm is tailored for Test Case Prioritisation. This algorithm was presented by Marijan, Gotlieb and Sen [24] as part of a case study to improve the testing efficiency of industrial video conferencing software. Unlike the previous algorithms that only take code coverage into account, this algorithm also considers historical failure data and test execution time. The objective of this algorithm is twofold: select the test cases with the highest consecutive failure rate, whilst also maximising the number of executed test cases in a limited time frame. The below algorithm has been modified slightly, since the time frame is a domain-specific constraint for this particular industry case and irrelevant for this thesis. Since this is a prioritisation algorithm rather than a selection or minimisation algorithm, it yields a total ordering of all the test cases in the

test suite, ordered using a weighted function.

The modified version of the algorithm (pseudocode is provided in Algorithm 3) takes three inputs:

- the set of test cases to prioritise $TS = \{T_1, \dots, T_n\}$
- the execution time for each test case $E = \{E_1, \dots, E_n\}$
- the failure status for each test case over the previous m successive executions $F = \{F_1, \dots, F_n\}$, where $F_i = \{f_1, \dots, f_m\}$

The algorithm starts by creating an array P of length n , which contains the priority of each test case. The priority of each test case is initialised at zero. Next, an $m \times n$ failure matrix MF is constructed and filled using the following formula.

$$MF[i, j] = \begin{cases} 1 & \text{if test case } T_j \text{ passed in execution } (current - i) \\ -1 & \text{if test case } T_j \text{ failed in execution } (current - i) \end{cases}$$

This matrix MF is visualised in Table 2.1. This table contains the hypothetical failure rates of the last three executions of six test cases.

run	T_1	T_2	T_3	T_4	T_5	T_6
$current - 1$	1	1	1	1	-1	-1
$current - 2$	-1	1	-1	-1	1	-1
$current - 3$	1	1	-1	1	1	-1

Table 2.1: Visualisation of the failure matrix MF .

Afterwards, P is filled with the cumulative priority of each test case, which is calculated by multiplying its failure rate with a domain-specific weight heuristic ω . This heuristic is used to derive the probability of repeated failures of the same test, given earlier failures. In their paper [24], the authors apply the following weights:

$$\omega_i = \begin{cases} 0.7 & \text{if } i = 1 \\ 0.2 & \text{if } i = 2 \\ 0.1 & \text{if } i \geq 3 \end{cases}$$

$$P_j = \sum_{i=1 \dots m} MF[i, j] * \omega_i$$

Finally, the algorithm groups test cases based on their calculated priority in P . Every test case that belongs to the same group is equally relevant for execution in the current test run. However, within every test group the tests will differ in execution time E . The

final step is to reorder test cases that belong to the same group in such a way that test cases with a shorter duration are executed earlier in the group.

Algorithm 3 ROCKET algorithm

```

1: Input: Set  $TS = \{T_1, \dots, T_n\}$  of all test cases,
      Execution time  $E$  of every test case,
      Failure status  $FS$  for each test case over the previous  $m$  successive iterations.
2: Output: Priority of test cases  $P$ .
3:  $P \leftarrow \text{array}[1 \dots n]$  ▷ initially 0
4:  $MF \leftarrow \text{array}[1 \dots m]$ 
5: for all  $i \in 1 \dots m$  do
6:    $MF[i] \leftarrow \text{array}[1 \dots n]$ 
7:   for all  $j \in 1 \dots n$  do
8:     if test case  $T_j$  failed in run ( $current - i$ ) then  $MF[i][j] \leftarrow -1$ 
9:     else  $MF[i][j] \leftarrow 1$ 
10: for all  $j \in 1 \dots n$  do
11:   for all  $i \in 1 \dots m$  do
12:     if  $i = 1$  then  $P[j] \leftarrow P[j] + (MF[i][j] * 0.7)$ 
13:     else if  $i = 2$  then  $P[j] \leftarrow P[j] + (MF[i][j] * 0.2)$ 
14:     else  $P[j] \leftarrow P[j] + (MF[i][j] * 0.1)$ 
15:  $Q \leftarrow \{P[j] | j \in 1 \dots n\}$  ▷ distinct priorities
16:  $G \leftarrow \text{array}[1 \dots \text{Card}(Q)]$  ▷ initially empty sets
17: for all  $j \in 1 \dots n$  do
18:    $p \leftarrow P[j]$ 
19:    $G[p] \leftarrow G[p] \cup \{j\}$ 
20: Sort every group in  $G$  based on ascending execution time in  $E$ .
21: Sort  $P$  according to which group it belongs and its position within that group.

```

2.3 Adoption in testing frameworks

Some of the approaches discussed above have been integrated in existing software testing frameworks. This paper will now proceed by conducting an analysis of these frameworks and tools to analyse which optimisation features are available and how they were implemented.

2.3.1 Gradle and JUnit

Gradle¹ is a dependency manager and application framework for Java, Groovy and Kotlin projects. Gradle supports multiple plugins to automate tedious tasks, such as configuration management, testing and deploying. One of the supported testing integrations is JUnit², which is the most widely used unit testing framework by Java developers. JUnit 5 is the newest version which is still actively being developed as of today. The framework is integrated as the testing framework of choice in several other Java libraries and frameworks, such as Android and Spring. JUnit offers mediocre support for features that optimise the execution of the test cases, especially when used in conjunction with Gradle. The following three key features are available:

1. **Parallel test execution:** JUnit comes bundled with multiple test processors that are responsible for processing test classes and to execute the test cases. One of these test processors is the `MaxNParallelTestClassProcessor`, which is capable of running a configurable amount of test cases in parallel. This results in a major speed-up of the overall test suite execution.
2. **Prioritise failed test cases:** Another test class processor which is provided by Gradle, is the `RunPreviousFailedFirstTestClassProcessor`. This processor will prioritise test cases that have failed in the previous run, similar to the idea of the ROCKET-algorithm (subsection 2.2.3), albeit without taking into account the duration of these test cases.
3. **Test order specification:** JUnit allows the user to specify the order in which test cases will be executed³. By default, a random yet deterministic order is used. The order can be manipulated by annotating the test class with the `@TestMethodOrder`-annotation, or by annotating individual test cases with the `@Order(int)`-annotation. This feature can only be used to alter the order of test cases within the same test class, it is not possible to perform inter-test class reordering. This feature could be used to sort test cases based on their execution time.

¹<https://gradle.org>

²<https://junit.org>

³<https://junit.org/junit5/docs/current/user-guide/#writing-tests-test-execution-order>



Figure 2.4: Logo of Gradle



Figure 2.5: Logo of JUnit 5

2.3.2 OpenClover

OpenClover⁴ is a code coverage framework for Java and Groovy projects. The framework was created by Atlassian and open-sourced in 2017. It profiles itself as “the most sophisticated code coverage tool”, by extracting useful metrics from the coverage results and by providing features that can optimise the test suite. Among these features are powerful integrations with development software and prominent Continuous Integration services. Furthermore, OpenClover offers the automatic analysis of the coverage results to detect relations between the application source code and the test cases. This feature allows OpenClover to predict which test cases will have been affected, given a set of modifications to the source code. Subsequently, these predictions can be interpreted to implement Test Case Selection. This results in a reduced test suite execution time.



Figure 2.6: Logo of Atlassian Clover

⁴<https://openclover.org>

Chapter 3

Proposed framework: VeloCity [TODO REVISE]

The implementation part of this thesis consists of a framework and a set of tools, tailored at optimising the test suite as well as providing accompanying metrics and insights. The framework was named *VeloCity* to reflect its purpose of enhancing the speed at which Continuous Integration is practised. This paper will now proceed by describing the design goals of the framework. Afterwards, a high-level schematic overview of the implemented architecture will be provided, followed by a more in-depth explanation of every pipeline step. In the final section of this chapter, the *Alpha* algorithm will be presented.

3.1 Design goals

VeloCity has been implemented with four design goals in mind:

1. **Extensibility:** It should be possible and straightforward to support additional Continuous Integration systems, programming languages and test frameworks. Subsequently, a clear interface should be provided to integrate additional prioritisation algorithms.
2. **Minimally invasive:** Integrating VeloCity into an existing test suite should not require drastic changes to any of the test cases.
3. **Language agnosticism:** This design goal is related to the framework being extensible. The implemented tools should not need to be aware of the programming language of the source code, nor the used test framework.
4. **Self-improvement:** The prioritisation framework supports all of the algorithms presented in section 2.2. It is possible that the performance of a given algorithm is strongly dependent on the nature of the project it is being applied to. In order to facilitate this behaviour, the framework should be able to measure the performance of every algorithm and “learn” which algorithm offers the best prediction, given a set of source code.

3.2 Architecture

The architecture of the VeloCity framework consists of seven steps that are performed sequentially in a pipeline fashion, as illustrated in the sequence diagram (Figure 3.1). Every step is executed by one of three individual components, which will now be introduced briefly.

3.2.1 Agent

The first component that will be discussed is the agent. This is the only component that depends actively on both the programming language, as well as the used test framework, since it must interact directly with the source code and test suite. For every programming language or test framework that needs to be supported, a different implementation of an agent must be provided. These implementations are however strongly related, so much code can be reused or even shared. In this thesis, an agent was implemented in Java, more specifically as a plugin for the widely used Gradle and JUnit test framework. This combination was previously described in subsection 2.3.1. This plugin is responsible for running the test suite in a certain prioritised order, which is obtained by communicating with the controller (subsection 3.2.2). After the test cases have been executed, the plugin sends a feedback report to the controller, where it is analysed.

3.2.2 Controller

The second component is the core of the framework, acting as an intermediary between the agent on the left side and the predictor (subsection 3.2.3) on the right side. In order to satisfy the second design goal and allow language agnosticism, the agent communicates with the controller using the HTTP protocol by exposing a *REST*-interface. Representational State Transfer [REST] is a software architecture used by modern web applications that allows standardised communication using existing HTTP methods. On the right side, the controller does not communicate directly with the predictor, but rather stores prediction requests in a shared database which is periodically polled by the predictor. Besides routing prediction requests from the agent to the predictor, the controller will also update the meta predictor by evaluating the accuracy of earlier predictions of this project.

3.2.3 Predictor and Metrics

The final component is twofold. Its main responsibility is to apply the prioritisation algorithms and predict an order in which the test cases should be executed. This order

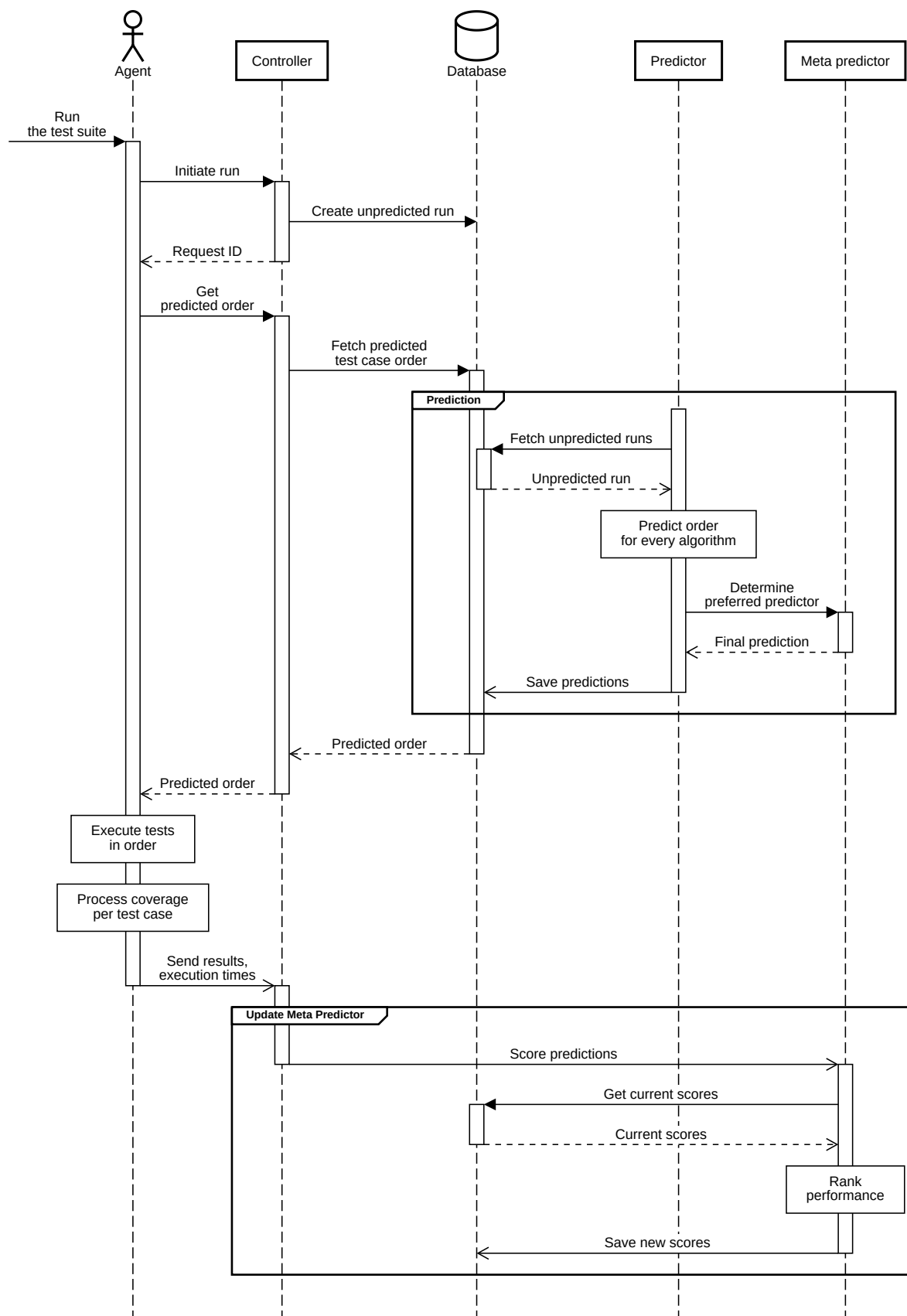


Figure 3.1: Sequence diagram of VeloCity

is calculated by first executing ten algorithms and subsequently picking the algorithm that has been preferred by the meta predictor. Additionally, this component is able to provide metrics about the test suite, such as identifying superfluous test cases by applying Test Suite Minimisation. More specifically, this redundancy is obtained using the greedy algorithm (subsection 2.2.1). Both of these scripts have been implemented in Python, because of its simplicity and existing libraries for many common operations, such as numerical calculations (NumPy¹) and machine learning (TensorFlow²).

3.3 Pipeline

This section will elaborate on the individual steps of the pipeline. The steps will be discussed by manually executing the pipeline that has hypothetically been implemented on a Java project. For the sake of simplicity, this explanation will assume a steady-state situation, ensuring the existence of at least one completed run of this project in the database at the controller side.

3.3.1 Initialisation

As was explained before, the provided Java implementation of the agent was designed to be used in conjunction with Gradle. In order to integrate VeloCity into a Gradle project, the build script (`build.gradle`) should be modified in two places. The first change is to include and apply the plugin in the header of the file. Afterwards, the plugin requires three properties to be configured:

- `base` the path to the Java source files, relative to the location of the build script. This will typically resemble `src/main/java`.
- `repository` the url to the git repository at which the project is hosted. This is required in subsequent steps of the pipeline, to detect which code lines have been changed in the commit currently being analysed.
- `server` the url at which the controller can be reached.

Listing 3.1 contains a minimal integration of the agent in a Gradle build script, applied to a library for generating random numbers³. The controller is hosted at the same host as the agent and is accessible at port 8080.

```
1 buildscript {  
2     dependencies {
```

¹<https://numpy.org/>

²<https://www.tensorflow.org/>

³<https://github.com/thepieterdc/random-java>

```
3      classpath 'io.github.thepieterdc.velocity:velocity-  
        junit:0.0.1-SNAPSHOT'  
4    }  
5  }  
6  
7  plugins {  
8      id 'java'  
9  }  
10  
11  apply plugin: 'velocity-junit'  
12  
13  velocity {  
14      base 'src/main/java/'  
15      repository 'https://github.com/thepieterdc/random-java'  
16      server 'http://localhost:8080'  
17  }
```

Listing 3.1: Minimal Gradle buildsript

After the project has been configured, the test suite must be executed. For the Gradle agent, this involves executing the built-in `test` task. This task requires an additional argument to be passed, which is the commit hash of the changeset to prioritise. In every discussed Continuous Integration system, this commit hash is available as an environment variable.

The first step is for the agent to initiate a new test run in the controller. This is accomplished by sending a `POST`-request to the `/runs` endpoint of the controller, which will reply with an identifier. On the controller side, this request will result in a new prioritisation request being enqueued in the database that will asynchronously be processed by the predictor daemon in the next step.

3.3.2 Prediction

The prediction of the test execution order is performed by the predictor daemon. This daemon continuously polls the database to fetch new test runs that need to be predicted. When a new test run is detected, the predictor executes every available prediction algorithm in order to obtain multiple prioritised test sequences. The following algorithms are available:

AllInOrder The first algorithm will simply prioritise every test case alphabetically and will be used for for benchmarking purposes in ??.

AllRandom The second algorithm has also been implemented for benchmarking purposes. This algorithm will “prioritise” every test case arbitrarily.

AffectedRandom This algorithm will only consider the test cases that cover source code lines which have been modified in the current commit. These test cases will be ordered randomly, followed by the other test cases in the test suite in no particular order.

GreedyCoverAll The first of three implementations of the Greedy algorithm (subsection 2.2.1) will execute the algorithm to prioritise the entire test suite.

GreedyCoverAffected As opposed to the previous greedy algorithm, the second Greedy algorithm will only consider test cases covering changed source code lines to be prioritised. After these test cases, the remaining test cases in the test suite will be ordered randomly.

GreedyTimeAll Instead of greedily attempting to cover as many lines of the source code using as few tests as possible, this implementation will attempt to execute as many tests as possible, as soon as possible. In other words, this algorithm will prioritise test cases based on their average execution time.

HGSAII This algorithm is an implementation of the algorithm presented by Harrold, Gupta and Soffa (subsection 2.2.2). It is executed for every test case in the test suite.

HGSAffected Similar to the *GreedyAffected* algorithm, this algorithm is identical to the previous *HGSAII* algorithm besides that it will only prioritise test cases covering changed source code lines.

ROCKET The penultimate algorithm is a straightforward implementation of the pseudocode provided in subsection 2.2.3.

Alpha The final algorithm has been inspired by the other implemented algorithms. section 3.4 will further elaborate on the details.

Subsequently, the final prioritisation order is determined by applying the meta predictor. Essentially, the meta predictor can be seen as a table which assigns a score to every algorithm, indicating its performance on this codebase. subsection 3.3.4 will explain later how this score is updated. The predicted order by the algorithm with the highest score is eventually elected by the meta predictor as the final prioritisation order, and saved to the database.

3.3.3 Test case execution

Regarding the agent, the identifier obtained in subsection 3.3.1 is used to poll the controller by sending a GET request to `/runs/id`, which will reply with the test execution order if this has already been determined. One of the discussed features of Gradle in subsection 2.3.1 was the possibility to execute test cases in a chosen order by adding annotations. However, this feature cannot be used to implement the Java agent, since it only supports ordering test cases within the same test class. In order to facilitate complete control over the order of execution, a custom `TestProcessor` and `TestListener` have been implemented.

The `TestProcessor` is responsible for processing every test class in the classpath and forward it along with configurable options to a delegate processor. The final processor in this chain will eventually perform the actual execution of the test class. Since the delegate processors that are built into Gradle will by default execute every method in the test class, the custom processor needs to work differently. The implemented agent will first store every received test class into a list and load the class to obtain all test cases in the class using reflection. After all classes have been processed, the processor will iterate over the prioritised order. For every test case t in the order, the delegate processor is called with a tuple of the corresponding test class and an options array which excludes every test case except t . This will effectively forward the same test class multiple times to the delegate processor, but each time with an option that restricts test execution to the prioritised test case, resulting in the desired behaviour.

Subsequently, the `TestListener` is a method that is called before and after every invocation of a test case. This listener allows the agent to calculate the duration of every test case, as well as collect the intermediary coverage and save this on a per-test case basis.

3.3.4 Post-processing and analysis

The final step of the pipeline is to provide feedback to the controller, to evaluate the accuracy of the predictions and thereby implementing the fourth design goal of self-improvement. After executing all test cases, the agent sends the test case results, the execution time and the coverage per test case to the controller by issuing a POST request to `/runs/id/test-results` and `/runs/id/coverage`.

Upon receiving this data, the controller will update the meta predictor using the following procedure. The meta predictor is only updated if at least one of the test cases has failed, since the objective of Test Case Prioritisation is to detect failures as fast as

possible, thus every prioritised order is equally good if there are no failures at all. If however a test case did fail, the predicted orders are inspected to calculate the duration until the first failed test case for every order. Subsequently, the average of all these durations is calculated. Finally, the score of every algorithm that predicted a below average duration until the first failure is increased, otherwise it is decreased. This will eventually lead to the most accurate algorithms being preferred in subsequent test runs.

3.4 Alpha algorithm

Besides the algorithms which have been presented in section 2.2, an additional algorithm has been implemented: the *Alpha* algorithm. This was constructed by examining the philosophy behind every discussed algorithm and subsequently combining the best ideas into a novel prioritisation algorithm. The specification below will assume the same naming convention as described in 5. The pseudocode is provided in Algorithm 4

The algorithm consumes the following inputs:

- the set of all n test cases: $TS = \{T_1, \dots, T_n\}$
- the set of m *affected* test cases: $AS = \{T_1, \dots, T_m\} \subseteq TS$. A test case t is considered “affected” if any source code line which is covered by t has been modified or removed in the commit that is being predicted.
- C : the set of all lines in the application source code, for which a test case $t \in TS$ exists that covers this line and that has not yet been prioritised. Initially, this set contains every covered source code line.
- the failure status of every test case, for every past execution out of k executions of that test case: $F = \{F_1, \dots, F_n\}$, where $F_i = \{f_1, \dots, f_k\}$. $F_{tj} = 1$ implies that test case t has failed in execution *current* - j .
- the execution time of test case $t \in TS$ for run $r \in [1 \dots k]$, in milliseconds: D_{tr} .
- for every test case $t \in TS$, the set TL_t is composed of all source code lines that are covered by test case t .

The first step of the algorithm is to determine the execution time E_t of every test case t . This execution time is calculated as the average of the durations of every successful (i.e.) execution of t , since a test case will be prematurely aborted upon the first failed assertion, which introduces bias in the duration timings. In case t has never been executed successfully, the average is computed over every execution of t .

$$E_t = \begin{cases} \overline{\{D_{ti} | i \in [1 \dots k], F_{ti} = 0\}} & \exists j \in [1 \dots k], F_{tj} = 0 \\ \overline{\{D_{ti} | i \in [1 \dots k]\}} & \text{otherwise} \end{cases}$$

Next, the algorithm executes every affected test case that has also failed at least once in its three previous executions. This reflects the behaviour of a developer attempting to resolve the bug that caused the test case to fail. Specifically executing *affected* failing test cases first is required in case multiple test cases are failing and the developer is resolving these one by one, an idea which was extracted from the ROCKET algorithm (subsection 2.2.3). In case there are multiple affected failing test cases, the test cases are prioritised by increasing execution time. After every selected test case, C is updated by subtracting the code lines that have been covered by at least one of these test cases.

Afterwards, the same operation is repeated for every failed but unaffected test case, likewise ordered by increasing execution time. Where the previous step helps developers to get fast feedback about whether or not the specific failing test case they were working on has been resolved, this step ensures that other failing test cases are not forgotten and are executed early in the run as well. Similar to the previous step, C is again updated after every prioritised test case.

Research (subsection 4.4.1) has indicated that on average, only a small fraction (10 % – 20 %) of all test runs will contain failed tests, resulting in the previous two steps not being executed at all. Therefore, the most time should be dedicated to executing test cases that cover affected code lines. More specifically, the next step of the algorithm executes every affected test case, sorted by decreasing cardinality of the intersection between C and the lines which are covered by the test case. Conforming to the prior two steps, C is also updated to reflect the selected test case. As a consequence of these updates, the cardinalities of these intersections change after every update, which will ultimately lead to affected tests not strictly requiring to be executed. This idea has been adopted from the Greedy algorithm subsection 2.2.1.

In the penultimate step, the previous operation is repeated in an identical fashion for the remaining test cases, similarly ordered by the cardinality of the intersection with the remaining uncovered lines in C .

Finally, the algorithm selects every test case which had not yet been prioritised. Notice that these test cases do not contribute to the test coverage, as every test case that would incur additional coverage would have been prioritised already in the pre-

vious step. Subsequently, these test cases are actually redundant and are therefore candidates for removal by Test Suite Minimisation. However, since this is a prioritisation algorithm, these tests will still be executed and prioritised by increasing execution time.

Algorithm 4 Alpha algorithm for Test Case Prioritisation

```

1: Input: Set  $TS = \{T_1, \dots, T_n\}$  of all test cases,
   Set  $AS = \{T_1, \dots, T_m\} \subseteq TS$  of affected test cases,
   Set  $C$  of all source code lines that are covered by any  $t \in TS$ ,
   Execution times  $D_{tr}$  of every test case  $t$ , over all  $k$  runs  $r$  of that test case,
   Failure status  $FS$  for each test case over the previous  $m$  successive iterations,
   Sets  $TL = \{TL_1, \dots, TL_n\}$  of all source code lines that are covered by test case
    $t \in TS$ .
2: Output: Ordered list  $P$  of  $n$  test cases and their priorities.
3:  $P \leftarrow \text{array}[1 \dots n]$  ▷ initially 0
4:  $i \leftarrow n$ 
5:  $FTS \leftarrow \{t \mid t \in TS \wedge (F[t][1] = 1 \vee F[t][2] = 1 \vee F[t][3] = 1)\}$ 
6:  $AFTS \leftarrow AS \cap FTS$ 
7: for all  $t \in AFTS$  do ▷ sorted by execution time in  $E$  (ascending)
8:    $C \leftarrow C \setminus TL[t]$ 
9:    $P[t] \leftarrow i$ 
10:   $i \leftarrow i - 1$ 
11:  $FTS \leftarrow FTS \setminus AFTS$ 
12: for all  $t \in FTS$  do ▷ sorted by execution time in  $E$  (ascending)
13:   $C \leftarrow C \setminus TL[t]$ 
14:   $P[t] \leftarrow i$ 
15:   $i \leftarrow i - 1$ 
16:  $AS \leftarrow AS \setminus FTS$ 
17: while  $AS \neq \emptyset$  do ▷ any element from  $AS$ 
18:   $t_{max} \leftarrow AS[1]$ 
19:   $tl_{max} \leftarrow \emptyset$ 
20:  for all  $t \in AS$  do
21:     $tl_{current} \leftarrow C \cap TL_t$ 
22:    if  $|tl_{current}| > |tl_{max}|$  then
23:       $t_{max} \leftarrow t$ 
24:       $tl_{max} \leftarrow tl_{current}$ 
25:   $C \leftarrow C \setminus tl_{max}$ 
26:   $P[t_{max}] \leftarrow i$ 
27:   $i \leftarrow i - 1$ 
28:  $TS \leftarrow TS \setminus (AS \cup FTS)$ 
29: while  $TS \neq \emptyset$  do ▷ any element from  $TS$ 
30:   $t_{max} \leftarrow TS[1]$ 
31:   $tl_{max} \leftarrow \emptyset$ 
32:  for all  $t \in TS$  do
33:     $tl_{current} \leftarrow C \cap TL_t$ 
34:    if  $|tl_{current}| > |tl_{max}|$  then
35:       $t_{max} \leftarrow t$ 
36:       $tl_{max} \leftarrow tl_{current}$ 
37:   $C \leftarrow C \setminus tl_{max}$ 
38:   $P[t_{max}] \leftarrow i$ 
39:   $i \leftarrow i - 1$ 
return  $P$ 

```

Chapter 4

Evaluation

This chapter will evaluate the performance of the framework presented in the previous chapter. The first section introduces the two test subjects that will be used in subsequent experiments. The next section will restate the research questions formally and extend these. Afterwards, we will elaborate on the procedure of the data collection. The final section will provide answers to the research questions as well as present the results of applying Test Case Prioritisation to the test subjects.

4.1 Test subjects

4.1.1 Dodona

Dodona¹ is an open-source online learning environment created by Ghent University, which allows students from secondary schools and universities in Belgium and South-Korea to submit solutions to programming exercises and receive instant, automated feedback. The application is built on top of the Ruby-on-Rails web framework. To automate the testing process of the application, Dodona employs GitHub Actions (section 1.2.2) which executes the more than 450 test cases in the test suite and performs static code analysis afterwards. The application is tested using the default MiniTest testing framework and SimpleCov² is used to record the coverage of the test suite. Currently, the coverage ratio is approximately 89%. This analysis will consider builds between January 1 and May 17, 2020.

4.1.2 Stratego

The second test subject has been created for the Software Engineering Lab 2 course at Ghent University in 2018. The application was created for a Belgian gas transmission system operator and consists of two main components: a web frontend and a backend. This thesis will test the backend in particular since it is written in Java using the Spring framework. Furthermore, the application uses Gradle and JUnit to execute the 300 – 400 test cases in the test suite, allowing the Java agent (section 3.2.1) to be applied directly.

¹<https://dodona.ugent.be/>

²<https://github.com/colszowka/simplecov>

4.2 Research questions

We will answer the following research questions in the subsequent sections:

RQ1: What is the probability that a test run will contain at least one failed test case? The first research question will provide useful insights into whether a typical test run tends to fail or not. The expectancy is that the probability of failure will be rather low, indicating that it is not strictly necessary to execute every test case and therefore making a case for Test Suite Minimisation.

RQ2: What is the average duration of a test run? Measuring how long it takes to execute a typical test run is required to estimate the benefit of applying any form of test suite optimisation. We will only consider successful test runs, to reduce bias introduced by prematurely aborting the execution.

RQ3: Suppose that a test run has failed, what is the probability that the next run will fail as well? The ROCKET algorithm (section 2.2.3) relies on the assumption that if a test case has failed in a given test run, it is likely to fail in the subsequent run as well. This research question will investigate the correctness of this hypothesis.

RQ4: How can Test Case Prioritisation be applied to Dodona and what is the resulting performance benefit? This research question will investigate the possibility to apply the VeloCity framework to the Dodona project and analyse how quickly the available predictors can discover a failing test case.

RQ5: Can the Java agent be applied to Stratego? Since the testing framework used by Stratego should be supported natively by the Java agent, this research question will verify its compatibility. Furthermore, we will analyse the prediction performance, albeit with a small number of relevant test runs.

4.3 Data collection

4.3.1 Travis CI build data

We can answer the first three research questions by analysing data from projects hosted on Travis CI (section 1.2.2). This data has been obtained from two sources.

The first source comprises a database [8] of 35 793 144 log files of executed test runs, contributed by Durieux et al. The magnitude of the dataset (61.11 GiB) requires a big

data approach to parse these log files. Two straightforward MapReduce pipelines have been created using the Apache Spark³ engine, to provide an answer to the first and second research question respectively.

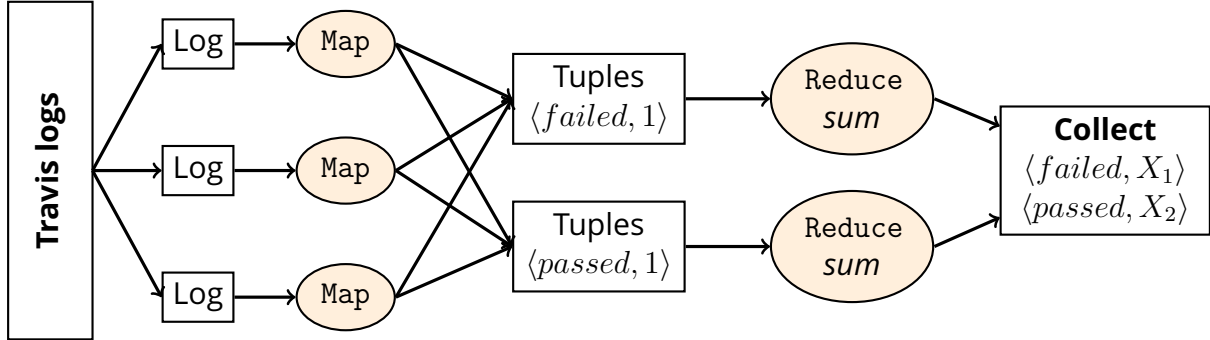


Figure 4.1: MapReduce pipeline to find the amount of failed test runs.

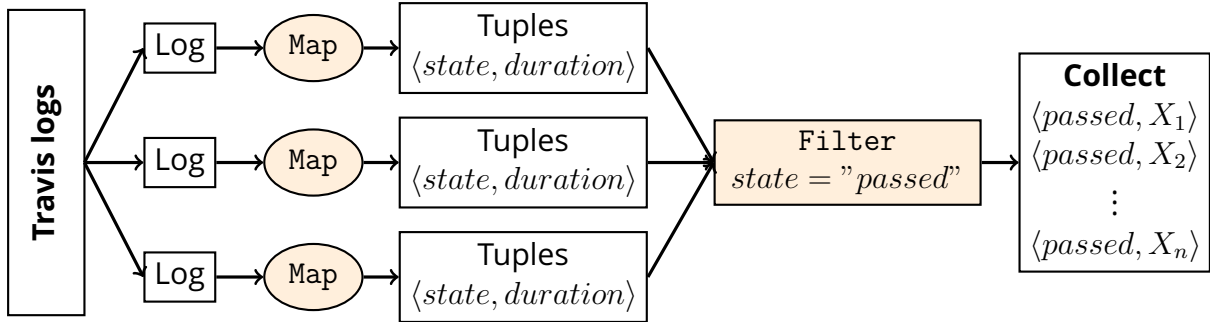


Figure 4.2: MapReduce pipeline to find the average duration of a successful test run.

In addition to the first source, another 3 702 595 jobs have been analysed from the *TravisTorrent* project [3]. To identify which projects are using Travis CI, the authors have crawled the GitHub API and examined the build status of every commit to retrieve the run identifier. Subsequently, the Travis CI API is used to obtain the build information, along with other useful statistics about the project. One of these additional values is the identifier of the previously executed run, which we can use to answer the second research question. Another interesting value is the amount of failed test cases. This value provides an accurate answer to the first research question since it indicates why the test run has failed. Without this information, the test suite might have failed to compile as opposed to an actual failure in the test cases. This dataset has been excluded from the second research question, as the included execution time does not correspond to the actual duration reported on the webpage of Travis CI. The authors have provided a Google BigQuery⁴ interface to allow querying the dataset more efficiently. Appendix A contains the corresponding executed queries.

³<https://spark.apache.org/>

⁴<https://bigquery.cloud.google.com/>

4.3.2 Dodona data

As mentioned before, Dodona utilises the MiniTest testing framework in conjunction with SimpleCov to calculate the coverage. MiniTest will by default only emit the name of every failed test case, without any further information. Furthermore, SimpleCov can only calculate the coverage for the entire test suite and does not allow us to retrieve the coverage on a per-test basis. To answer the fourth research question and apply the VeloClty predictors to Dodona, a Python script has been created to reconstruct the conditions of every failed test run. The script first queries the API of GitHub Actions to find which test runs have failed. This thesis will consider 120 failed runs. For every failed commit, the script retrieves the parent commit and calculates the coverage on a per-test basis. This thesis will assume that the coverage of the parent commit resembles the coverage of the failed commit. The coverage is calculated by applying the following two transformations to the parent commits and subsequently rescheduling these in GitHub Actions:

- **Cobertura formatter:** The current SimpleCov reports can only be generated as HTML reports, preventing convenient analysis. We can resolve this problem by using the Cobertura formatter instead, which generates XML reports. The controller already supports the structure of these reports, as this formatter is commonly used by Java testing frameworks as well.
- **Parallel execution:** The Dodona test suite currently executes the test cases by four processes concurrently, to reduce the execution time. Every process individually records the code coverage, and at the end of the test suite, SimpleCov merges these separate reports into one. However, this process is not entirely thread-safe since the test suite requires shared resources. We do not require thread-safety to calculate the total coverage, but we do require this to generate the coverage on a per-test basis. As a result, parallel execution has been disabled.

4.3.3 Stratego data

To integrate VeloClty with the existing Stratego codebase, we can use the instructions described in chapter 3. Afterwards, to analyse the prediction performance, we can take an approach similar to the previous test subject. The GitHub API has been used to identify the failed commits and to find their parent (successful) commits. The parent commits have subsequently been modified to use the VeloClty Java agent and have been executed using GitHub Actions.

4.4 Results

4.4.1 RQ1: Probability of failure

The two pie charts in Figure 4.3 illustrate the amount of failed and successful test runs. The leftmost chart visualises the failure rate in the dataset [8] by Durieux et al. 4 558 279 test runs out of the 28 882 003 total runs have failed, which corresponds to a failure probability of 18.74 %. The other pie chart uses data from the TravisTorrent [3] project. Since we can infer the cause of failure from this dataset, it is possible to obtain more accurate results. 42.89 % of the failed runs are due to a compilation failure where the test suite did not execute. For the remaining part of the runs, 225 766 out of 2 114 920 executions contain at least one failed test case, corresponding to a failure percentage of 10.67 %.

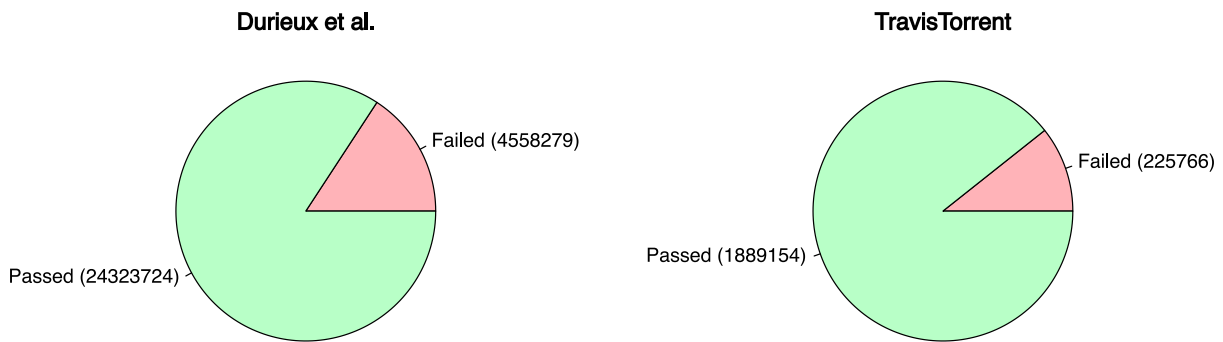


Figure 4.3: Probability of test run failure

4.4.2 RQ2: Average duration of a test run

The dataset by Durieux et al. [8] has been refined to only include test runs that did not finish within 10 s. A lower execution time generally indicates that the test suite did not execute and that a compilation failure has occurred instead. Table 4.1 contains the characteristics of the remaining 24 320 504 analysed test runs. The median and average execution times suggest that primarily small projects are Travis CI, yet the maximum value is very high. Figure 4.4 confirms that 71 378 test runs have taken longer than one hour to execute. Further investigation has revealed that these are typically projects which are using mutation testing, such as `plexus/yaks`⁵.

# runs	Minimum	Mean	Median	Maximum
24 320 504	10 s	385 s	178 s	26 h11 min26 s

Table 4.1: Characteristics of the test run durations in [8].

⁵A Ruby library for hypermedia (<https://github.com/plexus/yaks>).

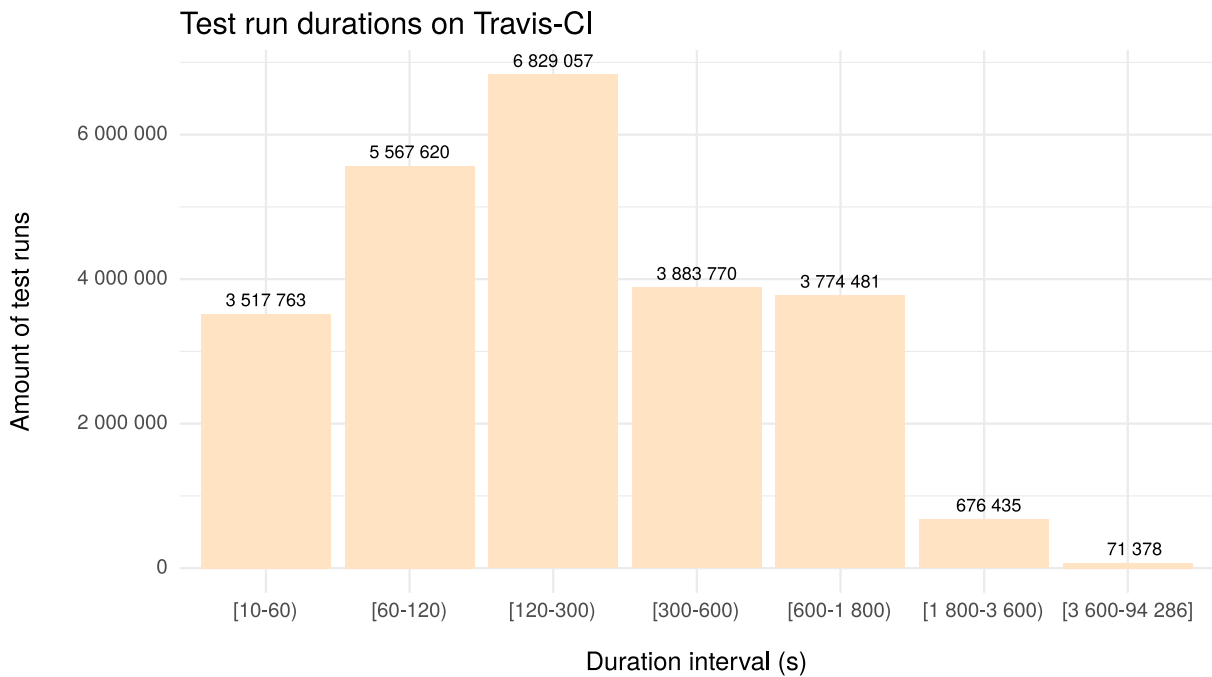


Figure 4.4: Test run durations on Travis CI

4.4.3 RQ3: Consecutive failure probability

Because the TravisTorrent project is the only dataset that contains the identifier of the previous run, only runs from this project have been used. This dataset consists of 211 040 test runs, immediately following a failed execution. As illustrated in Figure 4.5, 109 224 of these test runs have failed as well, versus 101 816 successful test runs (51.76 %).

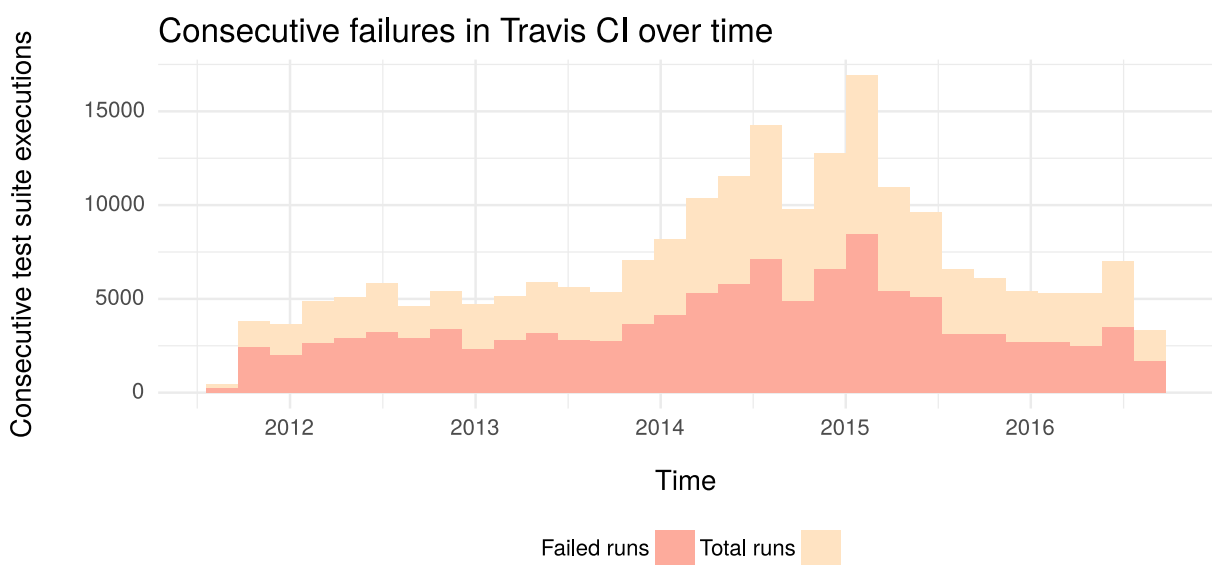


Figure 4.5: Consecutive test run failures on Travis CI

4.4.4 RQ4: Applying Test Case Prioritisation to Dodona

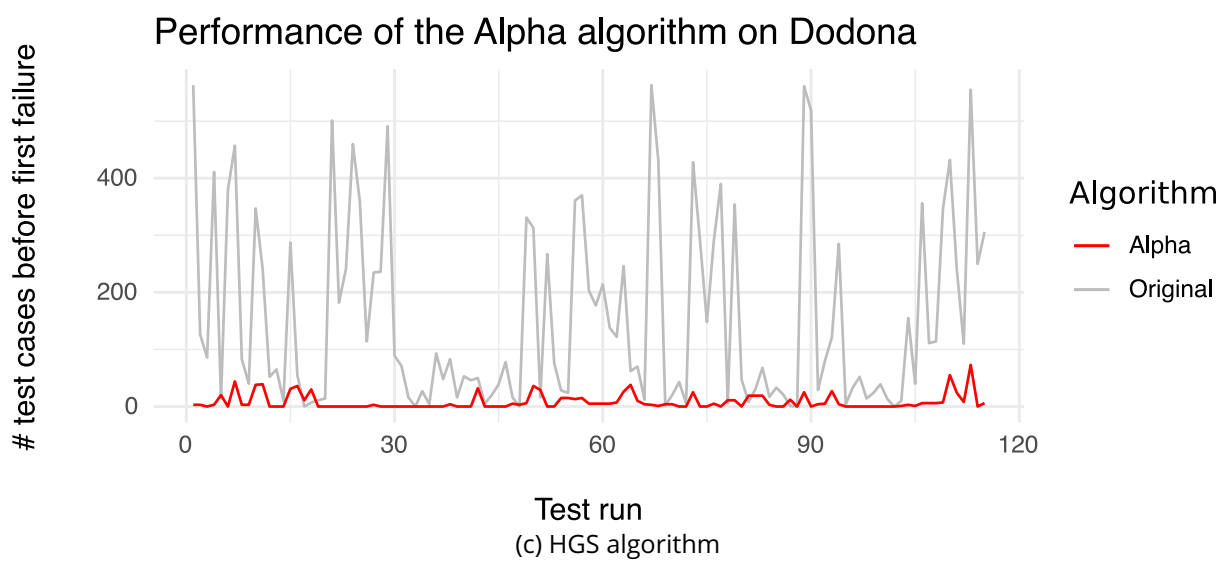
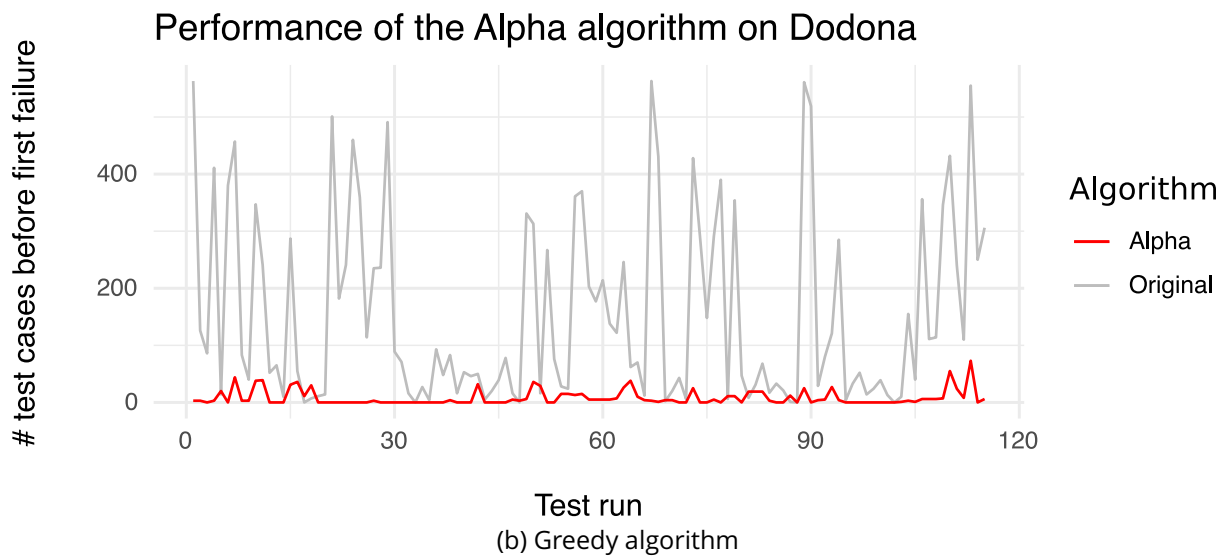
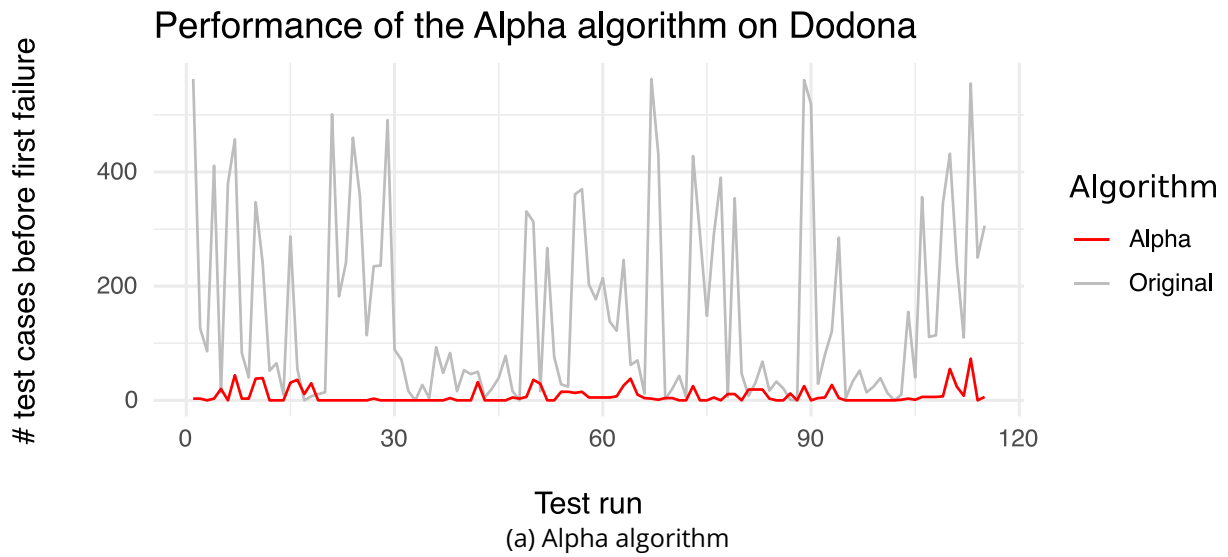
After executing the 120 failed test runs, the log files have been inspected. These log files have revealed that an error in the configuration was the actual culprit of five failed test runs, rather than a failed test case. These test runs have therefore been omitted from the results because the test suite did not execute. Since configuration-related problems require in-depth contextual information about the project, we cannot automatically predict these.

Table 4.2 contains the amount of executed test cases until we observe the first failure. These results indicate that every predictor is capable of performing at least one successful prediction. Furthermore, the maximum amount of executed test cases is lower than the original value, which means that every algorithm is a valid predictor. The data suggests that the Alpha algorithm and the HGS algorithm are the preferred predictors for the Dodona project. In contrast, the performance of the ROCKET algorithm is rather low.

Algorithm	Minimum	Mean	Median	Maximum
<i>Original</i>	0	155	78	563
Alpha	0	8	3	73
AffectedRandom	0	54	10	428
AllInOrder	0	119	82	460
AllRandom	0	90	27	473
GreedyCoverAffected	0	227	246	494
GreedyCoverAll	0	98	33	514
GreedyTimeAll	0	210	172	482
HGSAffected	0	61	10	511
HGSAll	0	124	54	507
ROCKET	0	210	170	482

Table 4.2: Amount of executed test cases until the first failure.

The previous results have been visualised in Figure 4.3. These charts confirm the low accuracy of the ROCKET algorithm. The Alpha algorithm and the HGS algorithm offer the most accurate predictions, with the former algorithm being the most consistent. Notice the chart of the Greedy algorithm, which succeeds in successfully predicting some of the test runs, while failing to predict others. This behaviour is specific to a greedy heuristic.



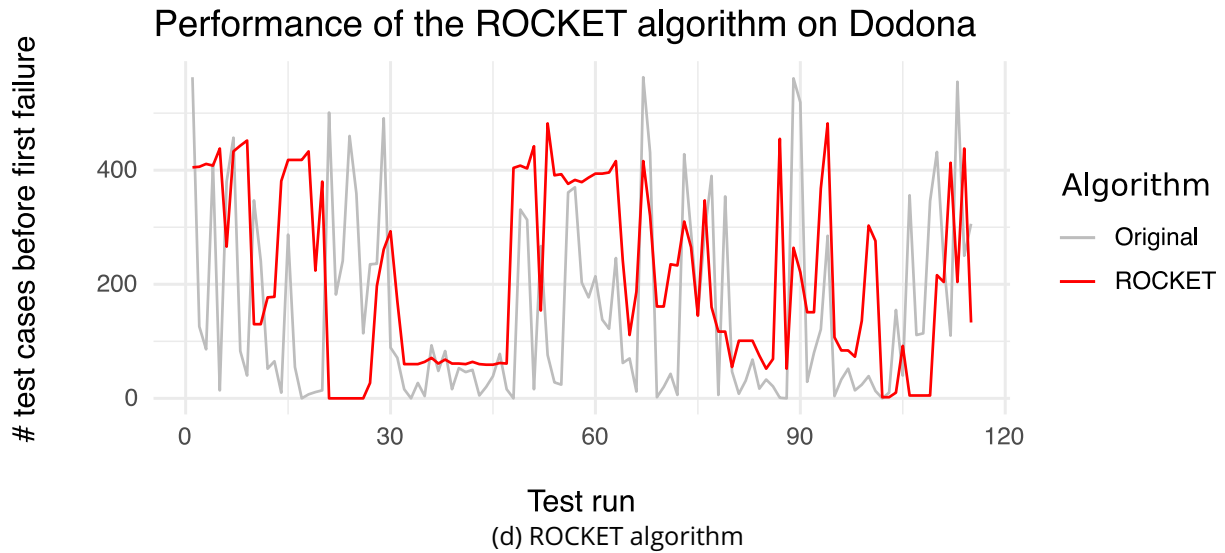


Figure 4.3: Prediction performance on the Dodona project.

The duration until the first observed failure is reported in Table 4.3. Observe that the previous table indicates that the ROCKET algorithm does not perform well, while this table suggests otherwise. We can explain this behaviour by examining the objective function of this algorithm. This function prioritises cases with a low execution time to be executed first.

Algorithm	Minimum	Mean	Median	Maximum
<i>Original</i>	0 s	135 s	123 s	380 s
Alpha	0 s	3 s	1 s	33 s
AffectedRandom	0 s	28 s	5 s	190 s
AllInOrder	0 s	82 s	71 s	270 s
AllRandom	0 s	43 s	11 s	270 s
GreedyCoverAffected	0 s	88 s	86 s	314 s
GreedyCoverAll	0 s	46 s	12 s	280 s
GreedyTimeAll	0 s	55 s	32 s	175 s
HGSAffected	0 s	35 s	6 s	356 s
HGSAII	0 s	75 s	34 s	377 s
ROCKET	0 s	54 s	32 s	175 s

Table 4.3: Duration until the first failure for the Dodona project.

4.4.5 RQ5: Integrate VeloCity with Stratego

The data collection phase has already proven that the Java agent is compatible with Stratego. Since VeloCity is not yet able to predict test cases which have been added in the current commit, we can only use 35 of the 54 failed test runs.

Similar to the previous test subject, Table 4.4 lists how many test cases have been executed before the first observed failure. The table only considers the four main algorithms, since the actual prediction performance was only secondary to this research question, and we have only analysed a small number of test runs. The results suggest that every algorithm except the ROCKET achieves a high prediction accuracy on this project.

Algorithm	Minimum	Mean	Median	Maximum
<i>Original</i>	0	68	2	278
Alpha	0	10	2	57
GreedyCoverAll	0	11	3	57
HGSAll	0	9	4	50
ROCKET	0	42	27	216

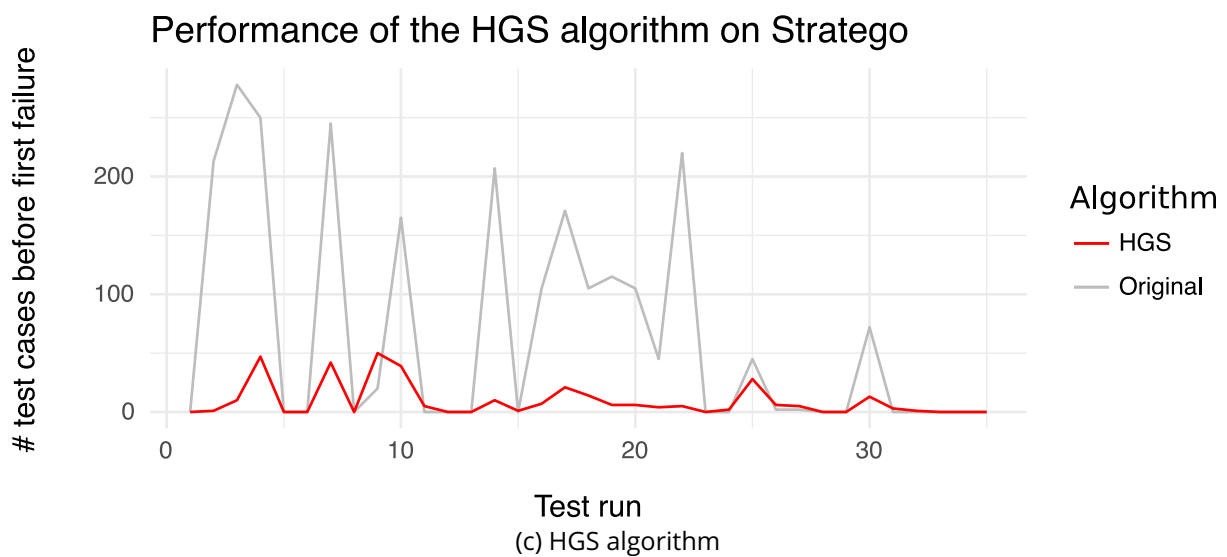
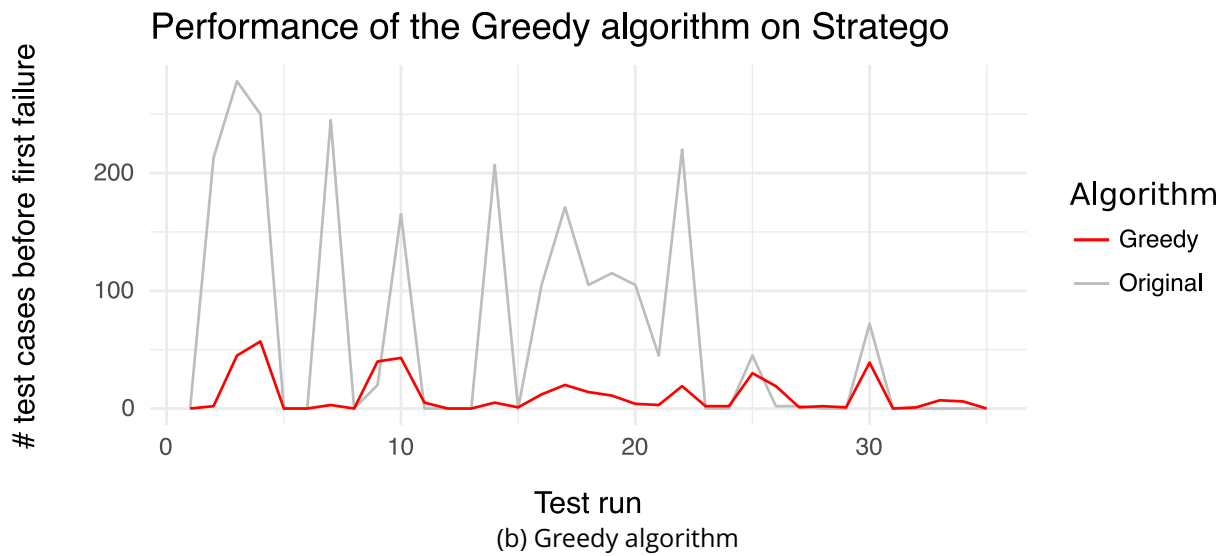
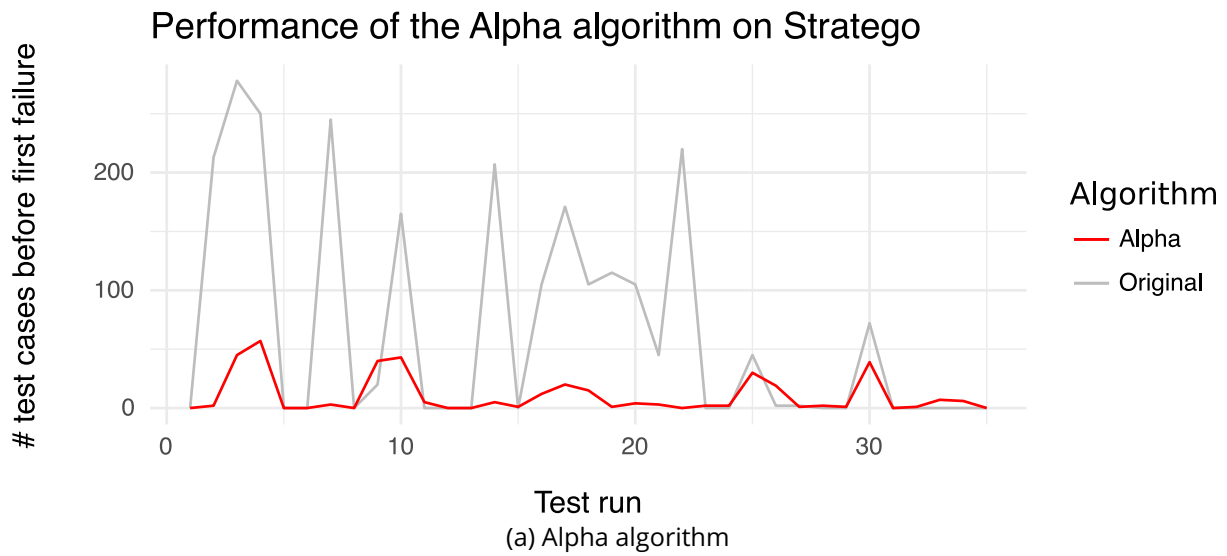
Table 4.4: Amount of executed test cases until the first failure.

Even though the performance of the ROCKET algorithm is suboptimal in the previous table, Table 4.5 does indicate that it outperforms every other algorithm time-wise. Notice that the predicted sequence of the HGS algorithm takes the longest to execute, while the previous table suggested a good prediction accuracy. The Alpha and Greedy algorithms seem very similar on both the amount of executed test cases, as well as the execution time.

Algorithm	Minimum	Mean	Median	Maximum
<i>Original</i>	0 s	62 s	8 s	233 s
Alpha	0 s	11 s	2 s	103 s
GreedyCoverAll	0 s	12 s	2 s	103 s
HGSAll	0 s	19 s	1 s	130 s
ROCKET	0 s	6 s	0 s	85 s

Table 4.5: Amount of executed test cases until the first failure.

Figure 4.1 further confirms the above statements. Notice the close resemblance of the charts of the Greedy algorithm and the Alpha algorithm, which indicates that a different failing test case is the cause of every test run failure. The ROCKET algorithm performs better on this project than on Dodona, yet not accurate.



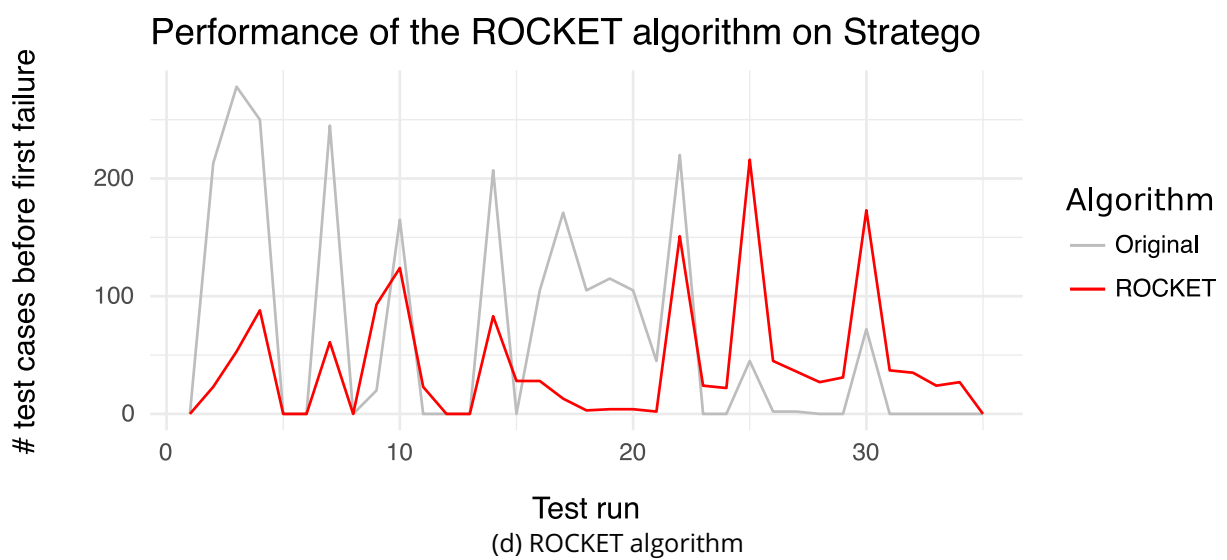


Figure 4.1: Prediction performance on the Stratego project.

Chapter 5

Conclusion

The main objective of this thesis was to study different approaches to optimise the test suite of a typical software project. Three approaches have been introduced to this extend: Test Suite Minimisation, Test Case Selection and Test Case Prioritisation. We have successfully implemented the latter approach using the VeloCity framework. Furthermore, this framework features the Alpha algorithm as a novel prioritisation algorithm. The performance of the Alpha algorithm has been evaluated, mainly on the UGent Dodona project. The results are promising, resulting in 95 % fewer executed test cases and 97 % less time spent waiting for the first test case to fail.

The second purpose of this thesis was to gain useful insights into the characteristics of a regular test suite, formalised into three research questions. The first question was to estimate the expected failure probability of a test run. To answer this question, we have analysed more than 28 million test runs on Travis-CI. This analysis has indicated that 18 % of those test runs have failed. Additionally, we have used this dataset to answer another question, which was to determine the typical duration of a test run. Statistical analysis has revealed that developers mainly use Travis-CI for small projects, with an average test suite duration of seven minutes. 0.20 % of the test suites take longer than one hour to execute, and some projects use mutation testing. The final question was to examine the probability of consecutive failing test runs. This probability was estimated at 52 % using a second Travis-CI dataset from the TravisTorrent project[3].

5.1 Future work

The proposed architecture currently features a Java agent, which supports the prediction of Gradle projects using ten available predictors. However, there is still room for improvements. The paragraphs below will suggest some ideas for possible enhancements.

5.1.1 Java agent

We can extend the functionality of the Java agent in multiple ways. Its current biggest weakness is the lack of support for parallel test case execution. To allow parallel testing, we must first solve a problem related to the scheduling process. Since the execution time of a test case can vary significantly, a coordination mechanism is required to schedule which test case should be executed on which thread. One possibility would be to consider the average execution time per test case, which we can obtain by examining prior runs. Alternatively, the scheduling can be performed at runtime using an existing inter-thread communication paradigm, such as message passing. Specific to the Java agent, implementing parallel execution requires us to modify the current `TestProcessor` to extend the `MaxNParallelTestClassProcessor` instead. A thread pool should ideally be used to diminish the overhead of restarting a new thread for every test case.

5.1.2 Predictions

We can make four different enhancements to the predictors.

For the first enhancement, the predictor should be able to discriminate between a unit test or an integration test. Recall that the scope of a unit test is limited to a small fraction of the application code and that its execution time is usually low. Contrarily, an integration test usually takes much longer to execute and tests multiple components of the application at once. The predictor should ideally make use of this distinction and assume that a failing unit test will almost certainly result in a failed integration test as well, and as such, prioritise unit tests over integration tests.

Secondly, the prediction algorithms currently take into account which source code lines have either been modified or removed to identify which test cases have been affected. Likewise, a change in the code of the test case itself should also consider that test case affected, as the change might have introduced a bug as well.

A third possible improvement would be to examine the performance of combining multiple prediction algorithms. Currently, the algorithms operate independently from each other, but there might be hidden potential in combining the individual strengths of these algorithms dynamically at runtime. A simple implementation is possible by modifying the existing meta predictor. Instead of assigning a score to the entire prediction, we could combine several predictions using predefined weights from earlier predictions.

Finally, the predictors do not currently consider branch coverage in addition to statement coverage. Not every coverage tool is capable of accurately reporting which branches have been covered, therefore this has not been implemented. Branch coverage can alternatively be supported by instrumenting the source code and rewriting every conditional expression as separate `if`-statements.

5.1.3 Meta predictor

The current implementation of the meta predictor increments the score of the predictor if the prediction was above-average, and decreases the score otherwise. However, a possible problem with this approach is that the nature of the source code might evolve and change as time progresses. As a result, it might take several test suite invocations for the meta predictor to prefer an alternative predictor. We can mitigate this effect if we would use a saturating counter (Figure 5.1) instead. This idea is also used in branch predictors of microprocessors and allows a more versatile meta predictor.

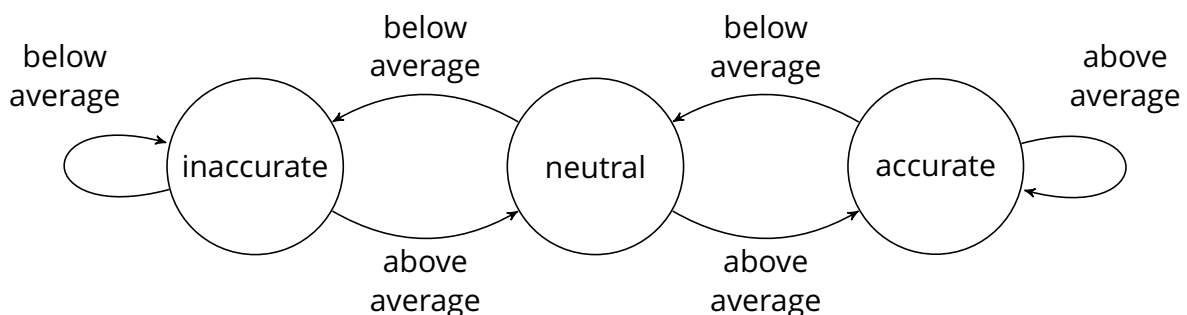


Figure 5.1: Saturating counter with three states

In addition to implementing a different update strategy, it might be worth to investigate the use of machine learning or linear programming models as a meta predictor, or even as a prediction algorithm.

5.1.4 Final enhancements

Finally, since we can apply every implemented algorithm to Test Suite Minimisation as well, we might extend the architecture to support this technique explicitly. Executing fewer test cases will result in even lower execution times.

Support for new programming languages and frameworks is possible by implementing a new agent. A naive implementation would be to restart the test suite after every executed test case, should test case reordering not be supported natively by the test framework.

Bibliography

- [1] *About GitHub Actions*. URL: <https://help.github.com/en/actions/getting-started-with-github-actions/about-github-actions>.
- [2] Mohammed Arefeen and Michael Schiller. "Continuous Integration Using Gitlab". In: *Undergraduate Research in Natural and Clinical Science and Technology (URN CST) Journal* 3 (Sept. 2019), pp. 1–6. DOI: 10.26685/urncst.152.
- [3] Moritz Beller, Georgios Gousios, and Andy Zaidman. "TravisTorrent: Synthesizing Travis CI and GitHub for Full-Stack Research on Continuous Integration". In: *Proceedings of the 14th working conference on mining software repositories*. 2017.
- [4] H.D. Benington. *Production of large computer programs*. ONR symposium report. Office of Naval Research, Department of the Navy, 1956, pp. 15–27. URL: <https://books.google.com/books?id=tLo6AQAAMAAJ>.
- [5] Thomas H. Cormen et al. *Introduction to Algorithms, Third Edition*. 3rd. The MIT Press, 2009. ISBN: 0262033844.
- [6] Michael Cusumano, Akindutire Michael, and Stanley Smith. "Beyond the waterfall : software development at Microsoft". In: (Feb. 1995).
- [7] Charles-Axel Dein. *dein.fr*. Sept. 2019. URL: <https://www.dein.fr/2019-09-06-test-coverage-only-matters-if-at-100-percent.html>.
- [8] Thomas Durieux et al. "An Analysis of 35+ Million Jobs of Travis CI". In: (2019). DOI: 10.1109/icsme.2019.00044. eprint: arXiv:1904.09416.
- [9] *Features • GitHub Actions*. URL: <https://github.com/features/actions>.
- [10] Michael R. Garey and David S. Johnson. *Computers and Intractability; A Guide to the Theory of NP-Completeness*. USA: W. H. Freeman & Co., 1990. ISBN: 0716710455.
- [11] *GitLab CI/CD*. URL: <https://docs.gitlab.com/ee/ci/>.
- [12] *GitLab Continuous Integration & Delivery*. URL: <https://about.gitlab.com/stages-devops-lifecycle/continuous-integration/>.
- [13] A. Govardhan. "A Comparison Between Five Models Of Software Engineering". In: *IJCSI International Journal of Computer Science Issues* 1694-0814 7 (Sept. 2010), pp. 94–101.
- [14] Standish Group et al. "CHAOS report 2015". In: *The Standish Group International* (2015). URL: https://www.standishgroup.com/sample_research_files/CHAOSReport2015-Final.pdf.

- [15] M. Jean Harrold, Rajiv Gupta, and Mary Lou Soffa. "A Methodology for Controlling the Size of a Test Suite". In: *ACM Trans. Softw. Eng. Methodol.* 2.3 (July 1993), pp. 270–285. ISSN: 1049-331X. DOI: 10.1145/152388.152391. URL: <https://doi.org/10.1145/152388.152391>.
- [16] Naftanaila Ionel. "AGILE SOFTWARE DEVELOPMENT METHODOLOGIES: AN OVERVIEW OF THE CURRENT STATE OF RESEARCH". In: *Annals of Faculty of Economics* 4 (May 2009), pp. 381–385.
- [17] "ISO/IEC/IEEE International Standard - Software and systems engineering –Software testing –Part 1:Concepts and definitions". In: *ISO/IEC/IEEE 29119-1:2013(E)* (Sept. 2013), pp. 1–64. DOI: 10.1109/IEEESTD.2013.6588537.
- [18] "ISO/IEC/IEEE International Standard - Systems and software engineering – System life cycle processes". In: *ISO/IEC/IEEE 15288 First edition 2015-05-15* (May 2015), pp. 1–118. DOI: 10.1109/IEEESTD.2015.7106435.
- [19] "ISO/IEC/IEEE International Standard - Systems and software engineering–Vocabulary". In: *ISO/IEC/IEEE 24765:2017(E)* (Aug. 2017), pp. 1–541. DOI: 10.1109/IEEESTD.2017.8016712.
- [20] Y. Jia and M. Harman. "An Analysis and Survey of the Development of Mutation Testing". In: *IEEE Transactions on Software Engineering* 37.5 (2011), pp. 649–678.
- [21] N. Landry. *Iterative and Agile Implementation Methodologies in Business Intelligence Software Development*. Lulu.com, 2011. ISBN: 9780557247585. URL: <https://books.google.be/books?id=bUHJAQAAQBAJ>.
- [22] G. Le Lann. "An analysis of the Ariane 5 flight 501 failure-a system engineering perspective". In: *Proceedings International Conference and Workshop on Engineering of Computer-Based Systems*. Mar. 1997, pp. 339–346. DOI: 10.1109/ECBS.1997.581900.
- [23] Simon Maple. *Development Tools in Java: 2016 Landscape*. July 2016. URL: <https://www.jrebel.com/blog/java-tools-and-technologies-2016>.
- [24] D. Marijan, A. Gotlieb, and S. Sen. "Test Case Prioritization for Continuous Regression Testing: An Industrial Case Study". In: *2013 IEEE International Conference on Software Maintenance*. 2013, pp. 540–543.
- [25] Robert C. Martin and Micah Martin. *Agile Principles, Patterns, and Practices in C# (Robert C. Martin)*. USA: Prentice Hall PTR, 2006. ISBN: 0131857258.
- [26] Bertrand Meyer. "Overview". In: *Agile!: The Good, the Hype and the Ugly*. Cham: Springer International Publishing, 2014, pp. 1–15. ISBN: 978-3-319-05155-0. DOI: 10.1007/978-3-319-05155-0_1. URL: https://doi.org/10.1007/978-3-319-05155-0_1.

- [27] Glenford J. Myers, Corey Sandler, and Tom Badgett. *The Art of Software Testing*. 3rd. Wiley Publishing, 2011. ISBN: 1118031962.
- [28] Dor Nir, Shmuel Tyszberowicz, and Amiram Yehudai. "Locating Regression Bugs". In: *Hardware and Software: Verification and Testing*. Ed. by Karen Yorav. Berlin, Heidelberg: Springer Berlin Heidelberg, 2008, pp. 218–234. ISBN: 978-3-540-77966-7.
- [29] Raphael Noemmer and Roman Haas. "An Evaluation of Test Suite Minimization Techniques". In: Dec. 2019, pp. 51–66. ISBN: 978-3-030-35509-8. DOI: 10.1007/978-3-030-35510-4_4.
- [30] A. Jefferson Offutt and Roland H. Untch. "Mutation 2000: Uniting the Orthogonal". In: *Mutation Testing for the New Century*. Ed. by W. Eric Wong. Boston, MA: Springer US, 2001, pp. 34–44. ISBN: 978-1-4757-5939-6. DOI: 10.1007/978-1-4757-5939-6_7. URL: https://doi.org/10.1007/978-1-4757-5939-6_7.
- [31] W. W. Royce. "Managing the Development of Large Software Systems: Concepts and Techniques". In: *Proceedings of the 9th International Conference on Software Engineering*. ICSE '87. Monterey, California, USA: IEEE Computer Society Press, 1987, pp. 328–338. ISBN: 0-89791-216-0. URL: <http://dl.acm.org/citation.cfm?id=41765.41801>.
- [32] John Ferguson Smart. *Jenkins: The Definitive Guide*. Beijing: O'Reilly, 2011. ISBN: 978-1-4493-0535-2. URL: <https://www.safaribooksonline.com/library/view/jenkins-the-definitive/9781449311155/>.
- [33] Travis. *Travis CI - Test and Deploy Your Code with Confidence*. Feb. 2020. URL: <https://travis-ci.org>.
- [34] Kristen R. Walcott et al. "TimeAware Test Suite Prioritization". In: *Proceedings of the 2006 International Symposium on Software Testing and Analysis*. ISSTA '06. Portland, Maine, USA: Association for Computing Machinery, 2006, pp. 1–12. ISBN: 1595932631. DOI: 10.1145/1146238.1146240. URL: <https://doi.org/10.1145/1146238.1146240>.
- [35] James Whittaker. "What is software testing? And why is it so hard?" In: *Software, IEEE* 17 (Feb. 2000), pp. 70–79. DOI: 10.1109/52.819971.
- [36] S. Yoo and M. Harman. "Regression Testing Minimization, Selection and Prioritization: A Survey". In: *Softw. Test. Verif. Reliab.* 22.2 (Mar. 2012), pp. 67–120. ISSN: 0960-0833. DOI: 10.1002/stv.430. URL: <https://doi.org/10.1002/stv.430>.

Appendices

Appendix A

TravisTorrent queries

```
1 SELECT  
2   COUNTIF(tr_log_bool_tests_failed) as failed ,  
3   COUNTIF(tr_log_bool_tests_ran) as ran,  
4   COUNT(1) as total  
5 FROM 'travistorrent'
```

Listing A.1: TravisTorrent query: Find the amount of failed runs

```
1 (  
2   SELECT gh_build_started_at , true as failed  
3   FROM 'travistorrent '  
4   WHERE  
5     tr_build_id IN (  
6       SELECT DISTINCT(tr_prev_build)  
7       FROM 'travistorrent '  
8       WHERE tr_log_bool_tests_ran=true AND  
9         tr_log_bool_tests_failed=true  
10    )  
11   AND gh_build_started_at IS NOT null AND  
12     tr_log_bool_tests_failed=true  
13 ) UNION ALL (  
14   SELECT gh_build_started_at , false as failed  
15   FROM 'travistorrent '  
16   WHERE tr_build_id IN (  
17     SELECT DISTINCT(tr_prev_build)  
18     FROM 'travistorrent '  
19     WHERE tr_log_bool_tests_ran=true AND tr_log_bool_tests_failed=  
20       false  
21   )  
22   AND gh_build_started_at IS NOT null AND tr_log_bool_tests_failed  
23     =true  
24 )
```

Listing A.2: TravisTorrent query: Find the probability of consecutive failures