

AUTOMATING TEST CASE IDENTIFICATION IN JAVA OPEN SOURCE PROJECTS ON GITHUB

Matej MADEJA, Jaroslav PORUBÄN, Michaela BAČÍKOVÁ
Matúš SULÍR, Ján JUHÁR, Sergej CHODAREV, Filip GURBÁL

Department of Computers and Informatics

Faculty of Electrical Engineering and Informatics

Technical University of Košice, Letná 9, 042 00 Košice, Slovakia

e-mail: {matej.madeja, jaroslav.poruban, michaela.bacikova,

matus.sulir, jan.juhar, sergej.chodarev, filip.gurbal}@tuke.sk

Abstract. Software testing is one of the very important Quality Assurance (QA) components. A lot of researchers deal with the testing process in terms of tester motivation and how tests should or should not be written. However, it is not known from the recommendations how the tests are written in real projects. In this paper, the following was investigated: (i) the denotation of the word “test” in different natural languages; (ii) whether the number of occurrences of the word “test” correlates with the number of test cases; and (iii) what testing frameworks are mostly used. The analysis was performed on 38 GitHub open source repositories thoroughly selected from the set of 4.3 M GitHub projects. We analyzed 20 340 test cases in 803 classes manually and 170 k classes using an automated approach. The results show that: (i) there exists a weak correlation ($r = 0.655$) between the number of occurrences of the word “test” and the number of test cases in a class; (ii) the proposed algorithm using static file analysis correctly detected 97% of test cases; (iii) 15% of the analyzed classes used `main()` function whose represent regular Java programs that test the production code without using any third-party framework. The identification of such tests is very complex due to implementation diversity. The results may be leveraged to more quickly identify and locate test cases in a repository, to understand practices in customized testing solutions, and to mine tests to improve program comprehension in the future.

Keywords: Program comprehension, Java testing, testing practices, test smells, open-source projects, GitHub

Mathematics Subject Classification 2010: 68-04

1 INTRODUCTION

The development of automated tests in a software project is a time-consuming and costly process, as it represents more than half of the entire development process [1]. The main aim of testing is to maintain the quality of the product and, in addition to that, tests describe the expected behavior of the production code being tested. Years ago, Demeyer et al. [2] suggested that if the tests are maintained together with the production code, their implementation is the most accurate mirror of the product specification and can be considered as up-to-date documentation. Tests can contain many useful production code metadata that can support program comprehension.

Understanding the code is one of the very first tasks a developer should cope with before the implementation of a particular feature. When the product specification changes (e.g., the requirements for new features are added), the developer must first understand them, then create his/her mental model [3] and finally, the created mental model is expressed in a specific artifact – code implementation. The problem is that two developers are likely to create two different mental models for the same issue because according to Mayer [4] mental model may vary with respect to its completeness and veridicality. A comprehension gap could arise when one developer needs to adapt another programmer's mental model from the code.

An assumption can be made that by using the knowledge about the structure and semantics of tests and their connection to the production code, it is possible to increase the effectiveness of program comprehension and reduce the comprehension gap. This would be possible, for example, by enriching the source code with metadata from the tests directly into the production code, e.g. data used for testing, test scenarios, objects relations, comments, etc. To achieve this goal, it is necessary to know in detail how the tests are actually written and what data they use.

There exist many guidelines on how tests should be created. First, naming conventions may aid the readability and comprehension of the code. According to the empirical study by Butler et al. [5], developers largely follow naming conventions. Our previous research [6] shows that there is a relation between the naming of identifiers in the test code and the production code being tested. This indicates that the relationship between the test and production code is not only at the level of method calls, object instances, or identifier references, but also at the vocabulary level, depending on the domain knowledge and mental model of a tester/developer.

Furthermore, many authors [7, 8, 9] define best practices to simplify the test with the benefit of a faster understanding of the testing code and the identification of test failure. Some guidelines lead to avoiding test smells [10] because as reported by recent studies [11, 12], their presence might not only negatively affect the comprehension of test suites but can also lead to test cases being less effective in finding bugs in the production code. All mentioned approaches are only recommendations but do not really express how the tests are written in real projects. That means we know how tests should be written, but we do not know how they

are written in practice. Many researchers have tried to clarify the motivation of writing tests [13, 14, 15], the impact of test-driven development (TDD) on code quality [16, 17] or the popularity of testing frameworks [18].

To reveal testing practices in real and independent projects it is necessary to find a way to identify test cases in a project, without the time-consuming code analysis. Much more important than the number of test cases is the information where they are located. When a testing framework is used, the test identification is mostly straightforward, e.g. by the presence of the framework imports. On the other hand, to obtain a general overview of testing practices regardless of the used framework, it is advisable to consider tests that do not use any third-party framework and can be regarded as customized testing solutions. In most of the related works, tests are identified by searching specific file and folder names, or some specific keywords. Considering that these keywords usually included the word “test” and based on the authors’ experience of Java test cases development, it can be assumed that there is a relation between the word “test” and the number of test cases in a file. That means searching for the “test” string could be beneficial for faster test case identification. Based on the previous reasoning, this paper defines the following hypothesis and research question:

H 1. There is a strong correlation ($r \notin (-0.8, 0.8)$) between the number of occurrences of the word “test” in the file content and the number of test cases.

RQ 1. How many testing classes are implemented as customized testing solutions without using any third party framework?

This paper is focused exclusively on unit testing and analyzes 38 projects that have been carefully selected (see Section 3.4.2) from all GitHub projects with Java as a primary language (most of the code written in Java). Section 2 presents the current state and found gaps in the research. In Section 3, the research method is described, containing an examination of whether it is appropriate to search for tests using the word “test” due to different natural languages of developers, an overview of known testing frameworks, and a proposed algorithm for static code analysis to automate the identification of test cases. Section 4 summarizes the results, threats to validity are mentioned in Section 5, and conclusions can be found in Section 6.

2 STATE OF THE ART

Many researchers examine software testing but we still know little about the structure and semantics of test code. This chapter summarizes the related work of software testing from various perspectives.

Learning about real testing practices is a constant research challenge. The goal of such research is mostly to find imperfections and risks, learn, and make recommendations on how to prevent them and how to streamline their development. Leitner and Bezemer [19] studied 111 Java-based projects from GitHub that contain

performance tests. Authors identify tests by searching for one or more terms in the test file name or for the presence of popular framework import, solely in the `src/test` project directory. Selected projects were subjected to manual analysis, in which they monitored several metrics. The most important result for this paper was the fact that 103 projects also included unit tests, usually following standardized best practices. On the other hand, the performance testing approach of the same projects often appears less extensive and less standardized. Another finding was that 58 projects (52%) mix performance tests freely with their functional test suite, i.e., performance tests are in the same package, or even the same test file, as functional tests. Six projects implemented tests as the usage examples. Using a similar approach [19], in our case by searching for the word “test” and searching for imports of testing frameworks in all project’s Java files, we would like to analyze unit tests, but with a careful selection from all GitHub projects at a specific time, resulting in more relevant projects used for analysis.

Code coverage, also known as test coverage, is a very popular method for evaluating project quality. Ellims et al. [20] investigated the usage of unit testing in practice in three projects that authors evaluated as well-tested. Statement coverage was found to be indeed a poor measure of test adequacy. According to the findings of Hemmati [21], basic criteria such as statement coverage are a very weak metric, detecting only 10% of the faults. A test case may cover a piece of code but miss its faults. According to Hilton et al. [22], coverage can be beneficial in the code review process if a smaller part of the project is evaluated. By reducing coverage to a single ratio of the whole project, much valuable information could be lost. Kochhar et al. [23] performed an analysis of 100 large open-source Java projects showing that 31% of the projects have coverage greater than 50% and only 8% are greater than 75%.

Many experiments try to express the quality of tests by testing “mutants” [24], i.e., by modifying a program in small ways to create artificial defects. According to Gopinath et al. [25] mutants do not necessarily represent real bugs, therefore, they are not able to relevantly evaluate the quality of the test suite nor to find relations between the coverage and mutants’ reveal. However, there is a statistically significant correlation between code coverage and bug kill effectiveness of real software errors (non-mutants) [26]. The quality of the test suite is influenced by the way the mental model is expressed in the code, so examining real tests is more beneficial instead of using mutants.

The fact that unit tests are the most common test type in a project is confirmed by Cruz et al. [27]: 39% of 1 000 analyzed Android projects used unit tests. Another finding was that frequently updated projects were more aware of the importance of using automated tests than those updated several years ago. The adoption of tests has increased over the last few years, so focusing on information mining from the tests makes sense.

Another type of research was done by Munaiah et al. [28], who focused on the assessment of GitHub projects. They proposed a tool that can be used to identify repositories containing real engineered software projects. The aim was to eliminate

the repository noise such as example projects, homework assignments, etc. One of the metrics they use for assessment is unit test occurrence in the project using test ratio (number of source lines of code in test files to the number of source lines of code in all source files) to quantify the extent of the unit testing effort. Package imports of *JUnit* and *TestNG* frameworks were searched to identify tests in the project. This method could be useful when looking for the occurrence of specific testing frameworks in the code.

3 METHOD

First of all, it is necessary to find suitable projects containing test cases. Thus, metadata of all GitHub open-source projects was obtained via GHTorrent [29] (Section 3.1) due to their high availability. GHTorrent collects projects' metadata from GitHub, one of the biggest project-sharing platform in the world. The experiment was limited to projects with Java as the primary language. Searching for testing frameworks' imports [30] or files containing the word "test" in the filename [19] are common test class identification techniques.

Because our main goal for the future is to improve production code comprehension from a particular test case, we go deeper in this study and try to identify specific test cases (not only test classes), therefore, it is necessary to consider whether the searching for the word "test" is appropriate. Keep in mind, that the aim is not to count the number of test cases in a project. Otherwise, we could run tests via an automated build tool (e.g. ant, maven, or gradle) and collect the number of tests. In that case, the issue is that building such open-source projects often fails [31] and we need to build every single project and run tests what is a time-consuming task. In this paper, we try to count and especially find the location of such test cases.

Since the testing process can also be denoted by other keywords (e.g. verify¹, examine, etc.), an in-depth analysis (Section 3.2) of testing process denotation in various foreign languages was performed, which showed that searching for the word "test" is suitable. Due to the limitations of the GitHub Search API, it was possible to search only one word across all Github Java projects.

As the framework is assumed to influence developer thinking and test case implementation, a list of 50 unit testing frameworks for Java (Section 3.3) has been created. Because the goal is to detect customized testing practices compared with framework-based ones in existing projects, it is not possible to use an automated method, and since it is not possible to manually analyze all GitHub projects, we need to select the most suitable ones. Based on the meaning of the word "test" we assume that there will be a correlation between the number of occurrences of the word "test" (in file content or filename) and the number of test cases.

¹ See Mockito `verify()` method used for soft assertions: <https://javadoc.io/static/org.mockito/mockito-core/3.11.2/org/mockito/verification/VerificationMode.html>

Therefore, three datasets were created using the searching GitHub API for (Section 3.4):

1. the word “test” in filename,
2. the word “test” in file content,
3. frameworks’ imports in file content (38 frameworks).

Every single project was searched as mentioned above, 4.3 million projects in total. It is possible to expect that the more occurrences of the word “test” in the project, the more test cases will be present in it and the more we will learn from it in the future. Therefore, projects with the highest occurrence of the word “test” (in file content or filename) or with the highest occurrence of a specific framework’s import were selected for manual analysis. By searching for “test” regardless of the framework, we were also able to analyze testing practices without using any third-party framework. Because GitHub contains many projects that are not relevant, e.g. testing, homework, or cloned projects, rules for searching relevant projects have been defined (Section 3.4.2), resulting in a set of projects used for manual and automated analysis. A script for automated analysis was created to partially automate the identification of test cases (see Section 3.5). All methodology details are described in the following sections.

3.1 Data Source

To provide conclusions that are as general as possible, it would be ideal to analyze all types of projects, i.e. proprietary and open source. Because of limited access to proprietary projects, this experiment is focused exclusively on open source projects. GitHub² has become one of the most popular web-based services to host both proprietary and mostly open-source projects, therefore, we can consider it a suitable source of projects. It provides an open Application Programming Interface (API)³ allowing one to work with all public projects (with small exceptions).

To avoid the latency of the official API, the GitHub Archive project⁴ stores public events from the GitHub timeline and publishes them via Google BigQuery. Downloading via Google BigQuery is charged, therefore, *GHTorrent* [29] was used instead, which provides a mirror of GitHub projects’ metadata. It monitors the GitHub public event timeline, retrieves contents and dependencies of every event, and requests GitHub API to store project data into the database. That includes general info about projects, commits, comments, users, etc. The study data mining started in May 2019, therefore, the last MySQL dump⁵ `mysql-2019-05-01` has been used.

² <https://github.com/>

³ <https://docs.github.com/en/rest>

⁴ <https://www.gharchive.org/>

⁵ <https://ghtorrent.org/downloads.html>

3.2 Denotation of the Word “test”

Leitner et al. [19] searched for tests only in `src/test` directory and test classes identified manually. However, the tests can be placed in any project’s directory (e.g. Android⁶ uses `src/androidTest`). Another approach is to search for “*test*” string in filenames as executed by Kochhar et al. [15] because they assumed that the tests would be exclusively in files containing the case-insensitive “*test*” string. As in the previous case, best practices lead the developer to use “test” in the file name, but it is not mandatory. For this reason, the most accurate should be searching for the word “test” in the file content. Of course, firstly it is necessary to consider whether the word “test” is the right one for searching. Therefore, the meaning of the word “test” using Google Translate⁷ was verified in 109 different languages (all available by Google) as follows:

1. From English to foreign language and back to English

Using this method the most frequent⁸ meanings of the word “test” in a foreign language were obtained. By translating them back to English we found out which foreign language translations correspond to the original word “test”.

2. From foreign language to English and back to foreign language

The opposite approach was used to find whether the string “test” has a meaning in a particular foreign language. The word was translated into English and all its meanings were verified against the available translation alternatives in the given language.

Multiple translations ensured that the correct meaning of the word in a particular language was understood. Using the first method it was found out that word sets related to the testing process of different foreign languages are mostly translated as “test” in English, see Figure 1. This means that when a foreign developer would like to express something related to testing (e.g. to write a test case), he/she will use mostly the word “test”. In this meaning, it is the first choice when searching test cases by a string. Occasionally occurred meaning outside of testing area, e.g., *essay*, *audition* or *flier*. Because such meanings occurred only infrequently, they can be omitted. There were also 14 languages which did not include the word “test” in their reverse translation at all, but its meaning was rather denoting *examination*, *check* or *quiz*.

A total of 44 languages used non-Latin charset. For these languages, the second approach did not make sense to use. For the remaining languages, the meaning was completely identical in 43 languages and the same or similar meaning in 20 cases.

⁶ <https://developer.android.com/>

⁷ <https://translate.google.com/>

⁸ Frequency determined by Google Translate service, indicates how often a translation appears in public documents: 3 – high; 2 – middle; 1 – low frequency.

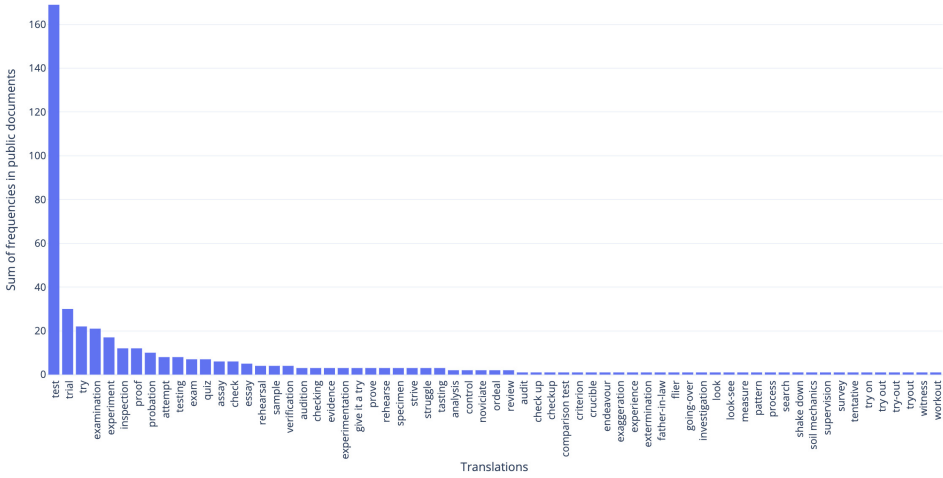


Figure 1. Sum of reverse translation frequency of the word “test” in public documents of different languages

We found only 2 languages (Hungarian⁹ and Latvian¹⁰), in which the word “test” has a completely different meaning, such as *body*, *hew*, or *tool* (nothing related to testing). The analysis shows that the word “test” will refer to the testing process in the code and the meaning can vary in very rare cases. Only the word “test” will be searched for in this study because of the rate limitations of the GitHub API (explained in Section 3.4).

3.3 Java Testing Frameworks

The crucial question is whether developers are motivated to use the word “test” in their code. The developer is often influenced by a testing framework, which leads him or her to different habits. As a part of this study, we analyzed 50 Java unit testing frameworks, extensions, and support libraries (see Table 1) to determine whether the use of the word “test” during test implementation is optional, recommended, or mandatory. The list was created from different sources, such as blogs, technical reports, research papers, etc.

Because it is sometimes difficult to find the boundary between unit and integration testing, the table lists frameworks supporting integration testing under the *unit testing* category. Information about the first version and the last commit may be interesting in terms of the framework lifetime and its occurrence in projects. Projects marked as *archived* or *test generators* in Table 1 were excluded from further analysis for the following reasons:

⁹ <https://translate.google.com/?sl=hu&tl=en&text=test>
¹⁰ <https://translate.google.com/?sl=lv&tl=en&text=test>

1. archived projects usually had unavailable documentation or were never released;
2. test generators produce tests that are not based on the programmer’s mental model but are generated automatically (semi-randomly), which is not interesting from the code comprehension point of view.

Name	Package for Import	Framework Type	First Version	Last Commit	Must Include "test"
SpryTest	N/A	U	N/A	N/A	N/A
Instinct	N/A	B	24.01.2007	07.03.2010 (archived)	N/A
Java Server-Side Testing framework (JSST)	N/A	U	17.11.2010	17.11.2010 (archived)	■
NUTester	N/A	U	05.02.2009	27.03.2012 (archived)	N/A
SureAssert	N/A	A	29.05.2011	04.02.2019 (archived)	N/A
Tacinga	N/A	U	14.02.2018	22.02.2018 (archived)	N/A
Unitils	N/A	U	29.09.2011	08.10.2015 (archived)	N/A
Cactus	org.apache.cactus	U	11.2008	05.08.2011 (archived)	■
Concutest	N/A	U	30.04.2009	12.01.2010 (archived)	■
Jtest	N/A	G	1997	21.05.2019 (last release)	■
Randoop	N/A	G	23.08.2010	05.05.2020	■
EvoSuite	N/A	G	25.12.2015 (v1.0.2)	30.04.2020	■
JWalk	N/A	G	19.05.2006	14.06.2017	■
TestNG	org.testng	U	31.07.2010 (v5.13)	11.04.2020	■
Artos	com.artos	U	22.09.2018	19.04.2020	■
JUnit 5	org.junit	U	10.09.2017	02.05.2020	■
JUnit 4	org.junit	U	16.02.2006	10.04.2020	■
JUnit 3	junit.framework	U	N/A	N/A	■
BeanTest	info.novatec.bean-test	U	23.04.2014	02.05.2015	■
GrandTestAuto	org.GrandTestAuto	U	21.11.2009	22.01.2014	■
Arquillian	org.jboss.arquillian	U	10.04.2012	21.04.2020	■
EtlUnit	org.bitbucket. bradleysmithllc.etlunit	U	02.12.2013 (v2.0.25)	04.04.2014	■
HavaRunner	com.github.havarunner	U	16.12.2013	08.06.2017	■
JExample	ch.unibe.jexample	U	2008	N/A	■
Cuppa	org.forgerock.cuppa	U	22.03.2016	01.10.2019	■
DbUnit	org.dbunit	U	27.02.2002	24.02.2020	■
GroboUtils	net.sourceforge.groboutils	U	20.12.2002	05.11.2004	■
JUnitEE	org.junit	U	23.07.2001 (v1.2)	11.12.2004	■
Needle	de.akquinet.jbosscc.needle	U	N/A	16.11.2016	■
OpenPojo	com.openpojo	U	13.10.2010	20.03.2020	■
Jukito	org.jukito	U/M	25.01.2011	17.04.2017	■
Spring testing	org.springframework.test	M/U	01.10.2002	06.05.2020	■
Concordion	org.concordion	U/SbE	23.11.2014 (v1.4.4)	27.04.2020	□
Jnario	org.jnario	B	23.07.2014		□
Cucumber-JVM	io.cucumber	B	27.03.2012	04.05.2020	□
Spock	spock.lang	B	05.03.2009	01.05.2020	□
JBehave	org.jbehave	B	2003	23.04.2020	□
JGiven	com.tngtech.jgiven	B	05.04.2014	10.04.2020	■

JDave	org.jdave	B	18.02.2008	17.01.2013	□
beanSpec	org.beanSpec	B	15.09.2007	27.06.2014	□
EasyMock	org.easymock.EasyMock	M	2001	10.04.2020	■
JMock	org.jmock	M	10.04.2007	23.04.2020	■
JMockit	org.jmockit	M	20.12.2012	13.04.2020	■
Mockito	org.mockito	M	2008	30.04.2020	■
Mockrunner	com.mockrunner	M	2003	16.03.2020	■
PowerMock	org.powermock	M	28.05.2014 (v1.5.5)	30.03.2020	■
AssertJ	org.assertj	A	26.03.2013	05.05.2020	■
Hamcrest	org.hamcrest	A	01.03.2012	06.05.2020	■
XMLUnit	org.xmlunit	A	03.2003	04.05.2020	■

Legend: U – unit; B – behavioural; A – assert; M – mock; G – generator; SbE – specification by example

Table 1: Analyzed unit testing frameworks and extensions for Java

It can be seen that 37 of 50 frameworks require the word “test” as method/class annotation (`@Test`) or part of its name (`testMethod`, `methodTest`). The listed frameworks are mostly extensions that depend on one of the base frameworks, such as *JUnit* or *TestNG*. Different versions of *JUnit* are listed separately because test labeling differs between them (annotations vs. method name format). A deeper analysis of frameworks’ JavaDocs revealed that many frameworks include other classes, methods, or annotations that include the word “test” in their names. Although the use of these methods is not mandatory, it may support the search.

3.4 Searching Projects and Data Gathering

The whole process of data gathering can be seen in Figure 2. GHTorrent provided 140 million GitHub projects. From this set all deleted, non-Java, or duplicated projects have been removed. After cleaning the initial data, a total of 6.7 million projects were kept for further analysis.

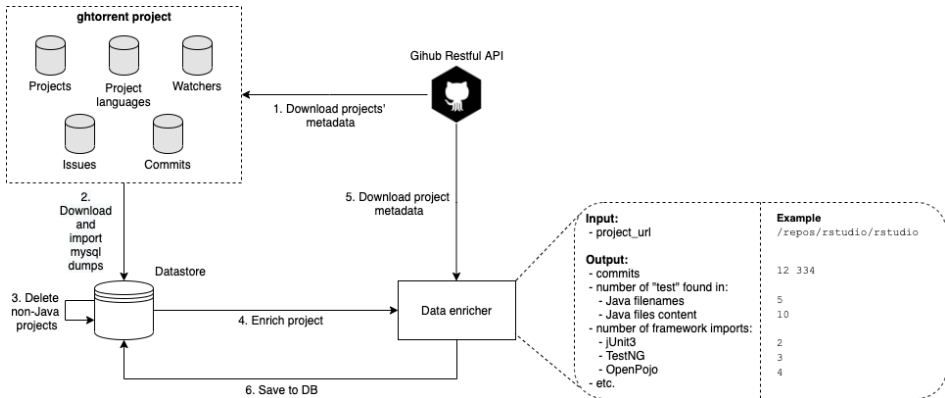


Figure 2. The GitHub data mining process for the study

GHTorrent contained only basic metadata about the projects, which was not sufficient for our needs. Given the meaning of the word “test” (see Section 3.2) we searched for it across all projects. The GitHub API provides a code search¹¹ endpoint, which indexes only original repositories. Repository forks are not searchable unless the fork has more stars than the parent repository. If the project has been detected as deleted, private, or blocked by GitHub during querying code search, it has been not considered. Finally, a total of 4.3 million projects were included. For each project, two requests to the GitHub code search API were called, as presented in Table 2. The GitHub code search API had the following limitations:

- up to 1 000 results for each search;
- up to 30 requests per minute (authenticated user);
- global requests rate limited at 5 000 requests per hour;
- only files smaller than 384 KB and repositories with fewer than 500 000 files are searchable.

Search “test” in	Example request at https://api.github.com/search/code
Java files content	?q=test+in:file+language:java+repo:apache/camel
Java filenames	?q=filename:test+language:java+repo:apache/camel

Table 2. The GitHub API requests used to search the string “test” in a project

3.4.1 Code Search Strategy

GitHub indexes only the default branch code (usually `master`), so the whole analysis was performed only using the default branch. The string “test” can also be a part of other words, e.g. *fastest*, *lastest*, *thisistestframework*. There exist 532 such words containing “test”¹² in total. To avoid inaccuracies when searching for a word of the selected string, false positives must be excluded from the search. When using regular GitHub search, the search term will appear in the results when driven by the following rules:

- string uses camel case convention without numbers¹³, e.g., `myTest`,
- string uses snake case convention, e.g., `my_test`, `test_123`;
- string includes a delimiter or special character (space, `.`, `#`, `$`, `@`, etc.), e.g., `test.delimiter`, `@Test`;
- search is case insensitive, e.g. `Test sentence`, `test sentence`.

¹¹ <https://docs.github.com/en/rest/reference/search>

¹² <https://www.thefreedictionary.com/words-containing-test>

¹³ Numbers can be used, but they are not considered as individual words, e.g. `123Test` or `test123` will not be found.

GitHub considers as Java language file any file with `.java` or `.properties` extensions. The same search rules apply to both search types: file content and filename search. Obviously, according to the above rules, GitHub search automatically filters the results, therefore, unwanted words containing the string “test” do not appear in the results, but neither the words `testing` or `testsAllMethods` will be matched.

3.4.2 Selection of Relevant Projects

When searching for different testing types, the effort is to go through as many projects as possible. Because GitHub contains millions of repositories, it is a challenge to choose the projects that can be the most instructive and filter out irrelevant ones. To make the selection as objective as possible, we planned to use *reaper* tool [28], which can assess a GitHub repository in collaboration with *GHTorrent* using project metadata and code: architecture, community, continuous integration, documentation, history, issues, license, and unit testing. By evaluating all these metrics (see [28] for details), *reaper* tags a particular repository as a real software project and thus exclude example projects, forks, irrelevant ones, etc.

Many assessment attributes of the *reaper* tool¹⁴ require project files to be available, so each project needs to be cloned or downloaded as an archive. For large projects, it can be gigabytes of data and the size of the project subsequently affects the length of the analysis. To find out whether *reaper* will be beneficial for our study, a manual analysis of 50 projects was performed and the results were compared with the evaluation by *reaper*. All available evaluation attributes were selected except for unit tests assessment because it was limited to *JUnit* and *TestNG* frameworks. The thresholds and weights of particular attributes defined by the developers of the tool were preserved because these values were considered empirically confirmed.

Because we want to select a sample of projects from which we would learn the most, projects with the highest number of files containing the word “test” in its body and filename were selected for the comparison. The same attributes as used by the *reaper* were taken into account in the manual evaluation, but the relevance of the project for this study was assessed by an observer. Evaluation of 50 projects using the *reaper* tool took 10 days, with the most time being spent on evaluating the project architecture. Many repositories with the highest “test” presence in file content or filename were actually identified as *Subversion* (SVN) mirrors¹⁵ by manual analysis and because there were multiple copies of the same code (caused by the SVN’s branching style), the projects were not relevant, but the *reaper* assessed such projects as suitable. According to this significant issue, important projects could be lost by assessing project in an automated manner, so it was concluded that

¹⁴ <https://github.com/RepoReapers/reaper>

¹⁵ e.g. <https://github.com/zg/jdk>, <https://github.com/dmatej/Glassfish>, <https://github.com/svn2github/cytoscape>

it is more efficient to select projects manually driven by the following rules, inspired by existing research:

- **Priority** was given to projects with the highest number of the word “test” in the project (in file content and filename). According to [32] we can expect the presence of tests in popular projects. If it is assumed that the word “test” will be correlated with the number of test cases in the project, large and long maintained projects are expected, which authors consider the best sample for the study.
- **History**, as evidence of sustained evolution. Projects under 50 commits were excluded (inspired by the *reaper*) because they represented small or irrelevant projects. Those projects that contained a large number of commits (more than 1 000 per day), considered committed by a robot, were also excluded.
- **Originality** was evaluated by comparing the `readme` file for similarities in other repositories. By such comparison, it is possible to detect clones and similar repositories [33]. Jiang et al. [34] found that developers clone repositories to submit pull requests, fix bugs, add new features, etc. The problem is that developers often do not create forks but project clones (a manual copy of a project), but `readme` file is often unchanged.
- **Community**, as evidence of collaboration, was assessed by the number of contributors in the project. The more developers participate in the project, the more likely it is that the (testing) code will be written in a different style.

3.4.3 Searching Java Testing Frameworks

We were inspired by the work of Stefan et al. [30], who searched for Java performance testing frameworks imports to assess performance testing practices. In our work we are interested in the impact of testing frameworks on test writing, so we also searched for imports of all testing frameworks in Table 1 (excluding generators and archived projects).

Using the search for imports we obtained projects with different testing frameworks. Only projects that contained the word “test” in the Java file body at least once were queried. Because there was a large number of requests (37 per single project), the project set was limited to 500 000, ordered by the number of Java files containing the word “test” in its body, using the following request:

```
https://api.github.com/search/code?q="org.testng"+in:  
file+language:java+repo:apache/camel
```

For each testing framework, we created a separate list of projects, sorted by the occurrence of the word “test” in the project, to find projects with a high number of test cases if possible. Original repositories of the searched framework were removed from the analysis (e.g. when searching for JUnit, the original JUnit framework repository was excluded). Subsequently, the selection of relevant projects was performed

according to the steps mentioned in the Section 3.4.2. For some frameworks, e.g. *JExample*¹⁶, which were created as a part of the research [35], no software repositories with business focus were found and as a consequence, it was necessary to include also example, homework, or cloned/forked ones, if the original one was not publicly available.

3.5 Repository Analysis

Three different data sets were received by searching via GitHub API:

1. the word “test” in filename,
2. the word “test” in file content,
3. frameworks’ imports in file content.

The first four relevant and top projects (highest “test” or framework’s import string occurrence) were manually investigated from each set to find out the test writing practices. The projects were cloned¹⁷ and to keep the consistency between the “test” search and the manual analysis, the project was reverted to the timestamp of GitHub API download using the following command:

```
git checkout `git rev-list -n 1
--before="<DOWNLOADED_AT>" "<DEFAULT_BRANCH>"`
```

For each project, all files with the word “test” in content or filename, or framework’s import in file content has been selected as possible option for manual analysis. The project files that contained the largest occurrence of the word “test” and framework’s import in their content (expected a higher number of tests) were analyzed first. During the investigation of tests from different authors and projects, we created an automated supportive method for detecting the number of test cases in a file. It does not require compiling the code, such as for computing code coverage, or building abstract syntax tree (AST), e.g. indexing in an IDE.

Regardless of the framework, it is advisable to investigate the count of the following attributes of a source file containing the word “test”:

Annotations @Test: very popular mostly thanks to *JUnit* and *TestNG*.

Methods containing test in the beginning of the name: best practices lead developers to use this convention (also for historical purposes).

Methods containing Test in the end of the name: an alternative of previous one.

Public methods: possibly all public methods of a test class can be considered as tests.

¹⁶ <https://github.com/akuhn/jexample>

¹⁷ `git clone`

Occurrence of main: customized testing solutions are executed via `main()`.

File path containing test: should relate to testing.

Classes containing \$ in the name: the character `$` in a class name mostly denotes a generated code¹⁸ that should not be analyzed.

Total number of test occurrence in file content: to reveal the relation between executable test cases and the word “test” presence in the content.

All listed metrics (counts of occurrence in a file) were saved for each analyzed file. The pseudocode for collecting mentioned metrics can be seen in Listing 1 (implementation available at GitHub¹⁹). The presented algorithm is partly the result of the study because it was created in parallel with the manual analysis. The manual analysis complements the algorithm implementation and vice versa. This algorithm was used to evaluate the test identification for each Java file containing the word “test”. Subsequently, the automated identification was checked during the manual analysis to determine the correct number of test cases and the metric used for the calculation (e.g., the number of annotations and public methods can be the same, but the relevant number of tests can only come from one of them). It is necessary to identify the number of particular test cases to link a specific test case with the unit under test (UUT) and its specific method. Each test case is likely to represent a unique use case and thus unique information to enrich the production code.

Algorithm `predictTests(filePath)`

Input: File path to analyze.

Output: List of statistical data

```

content := load filePath content and remove comments
nonClassContent := remove all class content, keep only content outside
    of it
    such as imports or class annotations
classContent := remove all content outside of the class block and keep
    only
    first-level methods without body using /\{([\{\}\}]+|(?R))*\}/

annotations := matches count of regex /@Test/ in classContent
startsWithTest := matches count of regex
    /public +.*void *.* +[Tt]est [a-zA-Z\\d$\\_]* *\\(/
    in classContent
endsWithTest := matches count of regex
    /public +.*void *.* +[a-zA-Z$\\_]{1}[a-zA-Z\\d$\\_]*Test *\\(/
    in classContent
publicMethods := matches count of regex /public +.*void +.*\\(/
    in classContent
includesMain := matches count of /public +static +void +main.*\\(/
    in classContent

hasDollar := if $ in filename, then true, else false
testInPath := if "/test" in filePath, then true, else false

```

¹⁸ <https://docs.oracle.com/javase/specs/jls/se11/html/jls-3.html#jls-3.8>

¹⁹ <https://github.com/madeja/unit-testing-practices-in-java/blob/master/AnalyzeProjectCommand.php>

```

if TestNG import found in content, then
  if @Test found in nonClassContent, then
    testCaseCount := publicMethods
  else
    testCaseCount := annotations
else if JUnit4 import found in content, then
  testCaseCount := annotations
else if JUnit3 import found in content, then
  testCaseCount := startsWithTest
else if startsWithTest > 0, then
  testCaseCount := startsWithTest
else if annotations > 0, then
  testCaseCount := annotations
else
  testCaseCount := 0

return annotations, startsWithTest, endsWithTest, publicMethods
      includesMain, hasDollar, testInPath, testCaseCount

```

Listing 1. Pseudocode of the algorithm for gathering metadata and identified number of tests in a Java source file

Gathered metadata about test case identification were analyzed from different perspectives. Test classes with the highest number of the following attributes were analyzed:

1. `@Test` annotations,
2. public methods with names starting with `test`,
3. public methods with names ending with `Test`,
4. `main` method,
5. word “test” occurrence.

For framework-dependent searches there was an additional analysis of files with the highest framework import occurrence in the content.

3.6 Hypothesis and Research Question Evaluation

Our null and alternative hypotheses are:

H_{null} 1 (H 1). There is **not** a strong correlation ($r \in (-0.8, 0.8)$) between the number of occurrences of the word “test” in the file content and the number of test cases in projects with high number of “test” occurrence.

H_{alt} 1 (H 1). There is a strong correlation ($r \notin (-0.8, 0.8)$) between the number of occurrences of the word “test” in the file content and the number of test cases.

The method of calculating standard Pearson’s correlation coefficient [36] was used to confirm or reject H 1. The correlation coefficient was calculated as follows:

$$r = \frac{\sum(x - m_x)(y - m_y)}{\sqrt{\sum(x - m_x)^2 \sum(y - m_y)^2}} \quad (1)$$

where m_x is the mean of the vector x (number of “test” occurrences in file) and m_y is the mean of the vector y (number of identified test cases in file). We will consider the H_{null} 1 as accepted when $r \in (-0.8, 0.8)$, as only absolute correlation higher than 0.8 is commonly considered significant.

To address RQ 1, a class/file will be considered a customized testing solution if the following conditions are met:

- Must include actual tests of production code.
- There is at least one occurrence of the word “test”.
- There is no framework import from Table 1.
- File contains `main()` function.

The conditions are based on Section 4.4.2 which shows that customized testing solutions were mostly implemented as common java programs using `main()` function without using any third party framework import.

4 RESULTS

Using the automated script all repositories’ files from Table 3 were processed, 38 repositories and 170 076 classes altogether, from which 803 classes and 20 340 test methods were manually investigated. Some special practices in terms of the structure of the testing code or the developer’s reasoning were observed. The first 4 projects from Table 3 represent repositories with the largest occurrences of the word “test” in the filename, another 4 in file content and other repositories represent the top import occurrence of a particular framework. The whole dataset of searching “test” via GitHub API can be found at Zenodo²⁰.

4.1 Accuracy of Automated Test Case Identification

To evaluate the precision of the algorithm from Listing 1, results were compared to manual test identification of 20 340 test cases across all three datasets. Accuracy of 95.72 % for test cases detection was achieved by automated identification considering only test methods, i.e., 95.72 % of all test cases were correctly identified. Considering all 28 975 methods of manually analyzed files (with non-testing ones) a total accuracy of 96.97 % was achieved with the sensitivity of

$$\text{Sensitivity} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}} = \frac{19\,600}{19\,600 + 62} = 0.9968 \quad (2)$$

and specificity of

$$\text{Specificity} = \frac{\text{true negatives}}{\text{true negatives} + \text{false positives}} = \frac{8\,498}{8\,498 + 815} = 0.9125. \quad (3)$$

²⁰ <https://doi.org/10.5281/zenodo.4566198>

Repository	Framework	Analyzed Classes		Analyzed Tests		Java KLOC	T_A
		A	M	A	M		
openjdk/client	testng, junit	30 410	130	30 410	1 661	5 149	20 798
SpoonLabs/astor	junit	30 331	36	30 331	1 548	2 338	13 324
apache/camel	junit	10 438	81	10 438	625	1 240	6 847
apache/netbeans	testng, junit	13 056	78	13 056	1 627	5 009	11 908
JetBrains/intellij-community	testng, junit	20 375	49	20 375	4 805	3 842	13 630
SpoonLabs/astor	testng, junit	30 331	44	30 331	5 883	2 338	13 324
corretto/corretto-8	testng, junit	13 688	10	13 688	1 659	3 638	10 792
aws/aws-sdk-java	junit	28 574	18	28 574	302	3 680	20 528
wildfly/wildfly	arquillian	5 109	24	5 109	123	548	3 553
eclipse-ee4j/cdi-tck	arquillian	4 758	30	4 758	139	97	2 748
resteasy/Resteasy	arquillian	2 821	13	2 821	144	220	1 675
keycloak/keycloak	arquillian	1 681	16	1 681	104	396	1 286
jsfunit/jsfunit	cactus	222	13	222	125	21	142
bleathem/mojarra	cactus	737	16	737	250	171	556
topcoder-platform /tc-website-master	cactus	1 635	8	1 635	42	366	1 199
apache/hadoop-hdfs	cactus	325	4	325	20	101	282
zanata/zanata-platform	dbunit	770	21	770	171	197	554
B3Partners/brmo	dbunit	145	18	145	37	47	106
gilbertoca/construtor	dbunit	145	18	145	64	24	53
sculptor/sculptor	dbunit	153	11	153	101	26	103
geotools/geotools	groboutils	3 424	5	3 424	5	1 272	3 659
notoriousre-i-d/ce-packager	groboutils	107	11	107	75	46	91
tliron/prudence	groboutils	16	2	16	3	13	11
MichaelKohler/P2	jexample	36	12	36	53	4	24
akuhn/codemap	jexample	132	15	132	286	41	112
wprogLK/TowerDefenceANTS	jexample	17	3	17	50	9	12
rbhama/Jboss-Files	needle	44	21	44	30	5	30
akquinet/mobile-blog	needle	19	10	19	33	2	10
s-case/s-case	needle	46	15	46	13	39	33
dbarton-uk/population-pie	needle	7	6	7	16	1	4
abarhub/rss	openpojo	26	2	26	3	6	20
BRUCELLA2 /Prescriptions-Scolaires	openpojo	25	19	25	40	10	18
jpmorganchase/tessera	openpojo	382	8	382	12	45	234
tensorics/tensorics-core	openpojo	161	3	161	1	24	85
orange-cloudfoundry /static-creds-broker	jgiven	21	11	21	33	2	16
eclipse/sw360	jgiven	175	4	175	51	56	161
Orchardir /FantasyWorldSimulation	jgiven	54	13	54	198	7	37
kodokojo/docker-image-manager	jgiven	11	5	11	8	3	8
Sum		170 076	803	363 730	20 340	31 033	127 973

Legend: A – processed automated; M – investigated manually; KLOC – kilo of lines of code;
 T_A – average time of automated test case detection in ms.

Table 3. Statistics of the investigated repositories

Most false positives and false negatives occurrences were caused by customized testing solutions, e.g., when tests were performed directly from the `main()` function by calling methods of the class. If the naming conventions of the called (testing) methods were not governed by the principles of frameworks (e.g., prepending method name with “*test*” or using public methods), not all test cases were detected in an automated way.

4.2 Correlation Between the Number of the Word “test” and the Number of Test Cases in a Class

The proposed algorithm was used to identify all tests in all Java classes of projects from Table 3. The script was used for all Java files that contained string “test” in the file content or the filename (in total 170 076 files). Figure 3 shows the correlation with the linear regression line of the word “test” and the number of test cases in a particular class. A standard Pearson’s correlation coefficient of $r = 0.655$ was reached (statistical significance $p = 0.0$, rounded on 5 decimal places), that means there is a weak correlation when considering absolute threshold $\alpha = 0.2$ defined in Section 3.6. Nevertheless, from the perspective of finding projects containing tests, this technique is beneficial and can help future experimenters to filter projects containing tests much faster. Because projects have different numbers of test classes and use different frameworks, the detailed ratio of the word “test” occurrence and test case presence per project can be found at GitHub²¹.

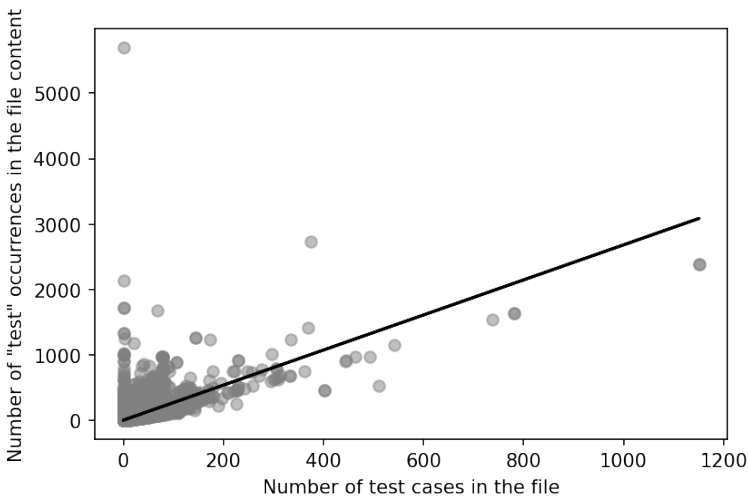


Figure 3. Correlation of the word “test” presence and number of test cases for analyzed classes by automated script

Due to existing research [19] that identified test files using searching “test” in the file path, when limiting our results to files containing “test” in the path (120 907 files) the correlation coefficient of $r = 0.6649$ was reached. On the other hand, 49 169 classes with 3 855 test cases were discarded. Limiting results to files containing “test” in filename (74 530 files), we reached correlation coefficient $r = 0.7004$ with

²¹ <https://github.com/madeja/unit-testing-practices-in-java/blob/master/correlation-boxplot.png>

loss of 95 546 classes and 17 440 test cases. By any limitation the correlation did not significantly change, therefore, to find as many test cases as possible it is convenient to search for the word “test” in the file content.

Occurrence of the function `main` without the third party testing framework (more explained in Section 4.4.2) was detected in 26 205 (15.41 %) classes containing the word “test” in their content. The proposed algorithm in Section 3.5 successfully identified test cases in only 6 % classes of this set. Because `main` tests make up a fairly large proportion and the identification of test cases is not clear, it is necessary to investigate this testing style deeper in the future.

H 1: There is a strong correlation ($r \notin (-0.8, 0.8)$) between the number of occurrences of the word “test” in the file content and the number of test cases.

We accept $H_{\text{null}} 1$ and reject $H_{\text{alt}} 1$ because only weak Pearson’s correlation coefficient $r = 0.655$ was achieved in general. In some projects, when the correlation was calculated for each project separately, a significant correlation was achieved but so far no relationship has been found concerning the framework, the number of the word “test” presence in the content, or other dependencies.

4.3 Efficiency of the Proposed Automated Test Case Identification

Executing a full code analysis, e.g. in an IDE, of a large project with thousands of kilo of lines of code (KLOC), is a time-consuming task. Such example is the project `openjdk/client` from Table 3. To get faster feedback about tests in a project, the proposed algorithm was used for static source code analysis. Because the proposed automated algorithm should run as a part of an integrated development environment (IDE) extension in the future it should be fast enough. To emulate a similar environment that a developer can use, a laptop with 2.3 GHz Dual-Core Intel Core i5 CPU and 8 GB 2133 MHz LPDDR3 RAM was used. In Table 3, the average time (T_A) of automated analysis executed 10 times can be seen. The average time of execution was 158ms per KLOC, which authors consider as a satisfactory response time in terms of user experience for use in an IDE extension.

4.4 Revealed Testing Practices

In related work (Section 2) there are best practices that developers can follow and therefore can be expected in the code. During the manual investigation of multiple repositories containing tests, we identified special testing practices used by developers, which are described in the following paragraphs. The listings that are given as examples come from the analyzed repositories, but the code was simpli-

fied for presentation purposes. Code listings refer to GitHub²² repository of this paper.

4.4.1 Testing Using Third Party Frameworks

Regular test. Tests that follow best practices and avoid test smells fall into this category. They represent the most of occurrences in the projects and since these approaches are already described in the available literature [7, 8, 9], this group will not be given detailed attention. However, the basic aspect of such tests is that information about context and evaluation are available directly in the particular test method (considering also test setup, teardown, and fixtures), thanks to which the test comprehension is straightforward.

Master test. This testing code style represents test classes which contain only one executable test method (see GitHub²³). *JUnit* will consider only the `all()` method as a test case because it is annotated with `@Test` annotation. Other methods are considered auxiliary ones. The problem with such a notation is the complexity of test comprehension. If the test fails, the developer only has information that the test case titled `all` failed but does not know what the test should have verified, what data was used, etc.

According to the best practices, it should be clear from the test name what the test verifies. In this context, from a semantic point of view, it is possible to consider methods as test cases on lines 1–8 (here from L1–8). The mentioned methods are crucial in terms of failure and understanding of the test, and from the method name, it is also clear what the test verifies. Another disadvantage of these test types is the *assertion roulette* test smell [10] because iterations of the test over the input data make it difficult to determine which data caused the test failure and whether the input data do not interfere with each other between the tests.

Reverse proxy test. If a separate test is written for each use case, the recommendations are met, but this does not mean that it will be easy to comprehend. Some tests call one auxiliary method in multiple tests and the result is evaluated in the auxiliary method. According to the test evaluation manner, they can be divided into:

1. Result evaluation via *method name* (see GitHub²⁴).
2. Result evaluation via *internal object state* (see GitHub²⁵).

²² <https://github.com/madeja/unit-testing-practices-in-java>

²³ https://github.com/madeja/unit-testing-practices-in-java/blob/master/examples/c_masterTest.java

²⁴ https://github.com/madeja/unit-testing-practices-in-java/blob/master/examples/c_reverseProxyMethod.java

²⁵ https://github.com/madeja/unit-testing-practices-in-java/blob/master/examples/c_reverseProxyObject.java

The first approach is much more difficult to comprehend due to the high degree of abstraction. It is not clear directly from the test method code (L6-8) what is compared during the test because the input data are loaded from a file determined by the test method name (L3). In the `JetBrains/intellij-community` project, from which the example is given, the `doTest()` method is the general one and it was necessary to investigate multiple classes to comprehend how tests are evaluated. At the same time, too generic auxiliary method can result in the *general fixture* test smell.

The second approach is similar to the previous one but uses the internal state of an object (that is initialized before a particular test during test setup) or the `enum` type with different method implementations. The problem may arise when object attribute or method input parameter change the control flow. If the same test is called with different object state or input data, the test logic does not change and therefore it is the same test. However, if the control flow changes in the test, e.g. by some variable value, it can be considered as a separate test (different flow, different test). If the same help method is called more than once, it may behave like two different test cases, which contradicts best practices and makes the comprehension difficult.

Multiple test execution. Server-side applications test different use cases, which require an action after the execution of base functionality, e.g. whether the right content is shown after main test execution (see GitHub²⁶). Because of using *JUnit3* in the example, every public method prepended by “test” is considered as test case, so `testEcho()` is executed twice; as a single test case and as a part of `testA4JRedirect()`.

4.4.2 Customized Testing Solutions

Custom testing practices are classic Java programs executable via `main()` function, whose task is to verify the functionality of the production code. Such tests are often written due to the possibility of configuring the execution via command line parameters, which allows variability of test execution. On the other hand, tests should not be so environmentally dependent that they need to be configured to such an extent. The second reason for writing such tests is that they make the code with a large number of test cases more readable. Test methods are called directly from `main()` and, if necessary, also the environment setup is performed in this function. The following ways of calling test methods and objects were observed (examples can be found at GitHub²⁷):

Calling methods one by one: all testing methods are manually called from `main()` together with parameters.

²⁶ https://github.com/madeja/unit-testing-practices-in-java/blob/master/examples/c_multipleExecution.java

²⁷ <https://github.com/madeja/unit-testing-practices-in-java>

Calling methods according to input data: by iterating the test data, specific tests are called based on the current data.

Helper function that returns an array of test cases: the helper method returns an array of instances created from abstract classes, whereas the abstract methods (which represent test cases) are implemented during the instance creation. The `main()` contains an iteration over the array of object instances.

Iterating values of enum: similar to the previous one, but it iterates over `enum` values. When creating the `enum`, the method of test class is implemented and the data is set. The test class has its own implementation of a method and state in each iteration.

Calling constructor: in the `main` function the testing class instance is created and the tests are called from the constructor.

There is a problem of how to identify such tests using an automated way and how to determine the number of tests in such a class. The `main()` function also occurs in classic tests (e.g. to run test outside of IDE or without a build automation tool²⁸), e.g. based on *JUnit* or *TestNG*. The function can also be found in modified runners of testing frameworks. To clearly distinguish the presence of a customized solution without any framework, it is possible to check the presence of the framework import – if a class contains the `main()` function and an import together, it is a runner or regular test based on the framework, not a customized solution.

Other interesting ways of writing customized tests were also observed. For example, in the `openjdk/client` repository, there were tests for trichotomous relations for which a custom `@Test` annotation was implemented (see GitHub²⁹). The annotation is used to indicate the test and, at the same time, to define the type of comparison in the method (L1, L4). Thanks to the word “test” usage, it is possible to detect the correct number of tests, in a similar way as for *JUnit*. In this example, the impact of third party framework on the developer’s customized solution is visible. There are many tests in the repository using standardized frameworks, therefore the usage of `@Test` annotation is a logical way of defining a test case. Writing tests manually using a framework would not be as effective and would be difficult to comprehend. On the other hand, such tests in large iterations can easily give rise to the *assertion roulette* test smell, which makes it difficult to identify a test failure.

While in the previous case the test was evaluated using asserts, some approaches have their own error handling. e.g. in the same repository for all *ResourceBundle* classes, a helper test class `RBTestFmwk` has been implemented, which represents a custom framework and test classes inherit from it. The framework provides the processing of the `main()` function parameters, performing tests, and processing results. The test methods to be performed are defined as input parameters. The

²⁸ <https://junit.org/junit4/faq.html>

²⁹ https://github.com/madeja/unit-testing-practices-in-java/blob/master/examples/c_main1.java

disadvantage is that when performing such tests, it is necessary to know the internal structure of the class, at least method names that need to be performed.

In general, the following risks were observed by analyzing other `main` testing methods:

Execution interruption: If a test fails, execution may be completely interrupted and no further tests will be performed (e.g. raised exception).

Failure identification: Because testing is often performed repeatedly over different data, it can be difficult to identify the exact cause of test failure and in some cases may require debugging the test code.

Dependence: Tests often use the same sources or data for testing and may affect the results of other tests. Also, the tests are often order-dependent and the test order randomness was not found in any repository.

Occurrence of the `main()` function without any third party testing framework was detected in 26 205 (15.41 %) classes containing the word “test” in their content. The proposed algorithm in Section 3.5 successfully identified test cases in only 6 % classes of this set. The set can contain not only testing code, but also a production one. Because such classes make up a fairly large proportion and the identification of test cases is not clear due to the high diversity of writing such tests, it is necessary to carry out an extensive study dealing solely with this issue, to find a way to precisely identify such test cases.

RQ 1: How many testing classes are implemented as customized testing solutions without using any third party framework?

A total of 15 % of classes were identified as customized testing solutions. The diversity of such tests is very high, therefore, future investigation is needed. This high incidence is probably caused by the nature of big projects with a high occurrence of the word “test” in file content and it is assumed the use of third party frameworks should be more common in smaller projects.

5 THREATS TO VALIDITY

Internal validity: The study relied on GHTorrent databank and GitHub API search algorithm to identify relevant projects. Because only projects with Java as a primary language were selected, testing practices in projects, where Java was not a major language could have been lost. Test classes that did not use the word “test” to indicate a test case were also lost. Searching for test cases was based on best practices and rules of the identified frameworks, but there may still exist other ways of how to identify them. The manual classification was based on observers’ experiences and identification of practices out of the generally known recommendations (best practices, test smells, etc.).

Test case detection results were compared to manual ones with an accuracy of 96.97%. As stated, it is necessary to further investigate customized testing solutions that use regular Java programs to test the production code. The implementation of such programs is often diametrically different and it is difficult to identify test cases. Real test cases were identified by the script in 6% of classes containing `main()` function.

External validity: To provide generalizable results, 20k of test cases were analyzed manually and 170k by an automated way. Also, the meaning and occurrence of the word “test” was analyzed for different natural languages and test frameworks. The results can be used to identify test cases in Java-based projects or projects with a different programming language with the usage of similar testing conventions. Despite the presented observations, our findings, as is usual in empirical software engineering, may not be directly generalized to other systems, particularly to commercial or to the ones implemented in other programming languages.

6 CONCLUSION AND FUTURE WORK

This paper presented an empirical study of Java open source GitHub projects to better understand how to identify test cases and their exact location in the project without the need for deep and time-consuming dynamic code analysis. An algorithm based on searching the word “test” in the repository files content or filenames was proposed and, at the same time, the unusual testing practices were investigated. In total 20340 test cases in 803 classes were investigated manually and 170k classes in an automated way. We summarise the most interesting findings from our study:

- There is not a strong correlation between the number of occurrences of the word “test” and the number of test cases in a class.
- Searching for the word “test” in the file content can be used to identify projects containing tests.
- Using static file analysis, the proposed algorithm can correctly detect 97% of test cases.
- Approximately 15% of the analyzed files contains “test” in the content together with `main()` function whose represent regular Java programs that test the production code without using any third-party framework. The success rate of identification of such test cases is very low because of implementation diversity.

Several test writing styles were found and classified, along with code samples of the analyzed repositories. Possible code comprehension defects were also mentioned. Based on these findings the following main contributions of this paper are concluded:

- Possibility of fast and automated test case identification together with the exact location in the project.

- Finding of correlation coefficient $r = 0.655$ between the number of occurrences of the word “test” and the number of test cases in a file, which allows to threshold projects or files for similar analysis.
- Overview of observed testing practices concerning the existence of customized testing solutions with an emphasis on places in testing code usable for mining information about the production code.

As future work, we plan to find a solution for accurate identification of test cases in customized solutions, probably based on training a machine learning model with manually labeled test cases of such testing solutions. We believe that studying testing practices can help comprehend the production code more easily. Gathered data could be used for training a machine learning model to automatically recognize tests by the structure and nature of the code. At the same time, we would like to focus on mining information from tests that could support the production source code comprehension and streamline the development process.

Acknowledgement

This work was supported by project VEGA No. 1/0762/19: Interactive Pattern-Driven Language Development.

REFERENCES

- [1] SCALABRINO, S.—LINARES-VÁSQUEZ, M.—POSHYVANYK, D.—OLIVETO, R.: Improving Code Readability Models with Textual Features. 2016 IEEE 24th International Conference on Program Comprehension (ICPC), 2016, pp. 1–10, doi: 10.1109/ICPC.2016.7503707.
- [2] DEMEYER, S.—DUCASSE, S.—NIERSTRASZ, O.: Object-Oriented Reengineering Patterns. Elsevier, 2002, doi: 10.1016/B978-1-55860-639-5.X5000-7.
- [3] CORRITORE, C. L.—WIEDENBECK, S.: Mental Representations of Expert Procedural and Object-Oriented Programmers in a Software Maintenance Task. International Journal of Human-Computer Studies, Vol. 50, 1999, No. 1, pp. 61–83, doi: 10.1006/ijhc.1998.0236.
- [4] MAYER, R. E.: The Psychology of How Novices Learn Computer Programming. ACM Computing Surveys, Vol. 13, 1981, No. 1, pp. 121–141, doi: 10.1145/356835.356841.
- [5] BUTLER, S.—WERMELINGER, M.—YU, Y.: Investigating Naming Convention Adherence in Java References. 2015 IEEE International Conference on Software Maintenance and Evolution (ICSME), 2015, pp. 41–50, doi: 10.1109/ICSM.2015.7332450.
- [6] MADEJA, M.—PORUBÄN, J.: Tracing Naming Semantics in Unit Tests of Popular Github Android Projects. In: Rodrigues, R., Janousek, J., Ferreira, L., Coheur, L., Batista, F., Oliveira, H. G. (Eds.): 8th Symposium on Languages, Applications and Technologies (SLATE 2019). OpenAccess Series in Informatics (OASIS), Vol. 74, 2019, pp. 3:1–3:13, doi: 10.4230/OASIS.SLATE.2019.3.

- [7] NAYYAR, A.: *Instant Approach to Software Testing: Principles, Applications, Techniques, and Practices*. BPB Publications, 2019.
- [8] LEWIS, W. E.: *Software Testing and Continuous Quality Improvement*. Second Edition. Auerbach Publications, 2004.
- [9] GARCIA, B.: *Mastering Software Testing with JUnit 5: Comprehensive Guide to Develop High Quality Java Applications*. Packt Publishing, 2017.
- [10] VAN DEURSEN, A.—MOONEN, L. M. F.—VAN DEN BERGH, A.—KOK, G.: Refactoring Test Code. Proceedings of the 2nd International Conference on Extreme Programming and Flexible Processes in Software Engineering (XP2001), 2001, pp. 92–95.
- [11] PERUMA, A.—ALMALKI, K.—NEWMAN, C. D.—MKAOUER, M. W.—OUNI, A.—PALOMBA, F.: On the Distribution of Test Smells in Open Source Android Applications: An Exploratory Study. Proceedings of the 29th Annual International Conference on Computer Science and Software Engineering (CASCON '19), 2019, pp. 193–202.
- [12] SPADINI, D.—PALOMBA, F.—ZAIMAN, A.—BRUNTINK, M.—BACCHELLI, A.: On the Relation of Test Smells to Software Code Quality. 2018 IEEE International Conference on Software Maintenance and Evolution (ICSME), 2018, pp. 1–12, doi: 10.1109/ICSME.2018.00010.
- [13] LINARES-VÁSQUEZ, M.—BERNAL-CARDENAS, C.—MORAN, K.—POSHYVANYK, D.: How Do Developers Test Android Applications? 2017 IEEE International Conference on Software Maintenance and Evolution (ICSME), 2017, pp. 613–622, doi: 10.1109/ICSME.2017.47.
- [14] BELLER, M.—GOUSIOS, G.—PANICHELLA, A.—ZAIMAN, A.: When, How, and Why Developers (Do Not) Test in Their IDEs. Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering (ESEC/FSE 2015), ACM, 2015, pp. 179–190, doi: 10.1145/2786805.2786843.
- [15] KOCHHAR, P. S.—THUNG, F.—NAGAPPAN, N.—ZIMMERMANN, T.—LO, D.: Understanding the Test Automation Culture of App Developers. 2015 IEEE 8th International Conference on Software Testing, Verification and Validation (ICST), 2015, pp. 1–10, doi: 10.1109/ICST.2015.7102609.
- [16] FUCCI, D.—ERDOGMUS, H.—TURHAN, H.—OIVO, M.—JURISTO, N.: A Dissection of the Test-Driven Development Process: Does It Really Matter to Test-First or to Test-Last? *IEEE Transactions on Software Engineering*, Vol. 43, 2017, No. 7, pp. 597–614, doi: 10.1109/TSE.2016.2616877.
- [17] BISSI, W.—SERRA SECA NETO, A. G.—PEREIRA EMER, M. C. F.: The Effects of Test Driven Development on Internal Quality, External Quality and Productivity: A Systematic Review. *Information and Software Technology*, Vol. 74, 2016, pp. 45–54, doi: 10.1016/j.infsof.2016.02.004.
- [18] ZEROUALI, A.—MENS, T.: Analyzing the Evolution of Testing Library Usage in Open Source Java Projects. 2017 IEEE 24th International Conference on Software Analysis, Evolution and Reengineering (SANER), 2017, pp. 417–421, doi: 10.1109/SANER.2017.7884645.
- [19] LEITNER, P.—BEZEMER, C.-P.: An Exploratory Study of the State of Practice of Performance Testing in Java-Based Open Source Projects. Proceedings of the 8th

- ACM/SPEC International Conference on Performance Engineering (ICPE '17), ACM, 2017, pp. 373–384, doi: 10.1145/3030207.3030213.
- [20] ELLIMS, M.—BRIDGES, J.—INCE, D.C.: Unit Testing in Practice. 15th International Symposium on Software Reliability Engineering, 2004, pp. 3–13, doi: 10.1109/ISSRE.2004.44.
- [21] HEMMATI, H.: How Effective Are Code Coverage Criteria? 2015 IEEE International Conference on Software Quality, Reliability and Security, 2015, pp. 151–156, doi: 10.1109/QRS.2015.30.
- [22] HILTON, M.—BELL, J.—MARINOV, D.: A Large-Scale Study of Test Coverage Evolution. Proceedings of the 33rd ACM/IEEE International Conference on Automated Software Engineering (ASE 2018), ACM, 2018, pp. 53–63, doi: 10.1145/3238147.3238183.
- [23] KOCHHAR, P.S.—LO, D.—LAWALL, J.—NAGAPPAN, N.: Code Coverage and Postrelease Defects: A Large-Scale Study on Open Source Projects. IEEE Transactions on Reliability, Vol. 66, 2017, No. 4, pp. 1213–1228, doi: 10.1109/TR.2017.2727062.
- [24] JUST, R.—JALALI, D.—INOZEMTSEVA, L.—ERNST, M.D.—HOLMES, R.—FRASER, G.: Are Mutants a Valid Substitute for Real Faults in Software Testing? Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering (FSE 2014), ACM, 2014, pp. 654–665, doi: 10.1145/2635868.2635929.
- [25] GOPINATH, R.—JENSEN, C.—GROCE, A.: Mutations: How Close Are They to Real Faults? 2014 IEEE 25th International Symposium on Software Reliability Engineering, 2014, pp. 189–200, doi: 10.1109/ISSRE.2014.40.
- [26] KOCHHAR, P.S.—THUNG, F.—LO, D.: Code Coverage and Test Suite Effectiveness: Empirical Study with Real Bugs in Large Systems. 2015 IEEE 22nd International Conference on Software Analysis, Evolution, and Reengineering (SANER), 2015, pp. 560–564, doi: 10.1109/SANER.2015.7081877.
- [27] CRUZ, L.—ABREU, R.—LO, D.: To the Attention of Mobile Software Developers: Guess What, Test Your App! Empirical Software Engineering, Vol. 24, 2019, No. 4, pp. 2438–2468, doi: 10.1007/s10664-019-09701-0.
- [28] MUNAIAH, N.—KROH, S.—CABREY, C.—NAGAPPAN, M.: Curating GitHub for Engineered Software Projects. Empirical Software Engineering, Vol. 22, 2017, pp. 3219–3253, doi: 10.1007/s10664-017-9512-6.
- [29] GOUSIOS, G.: The GHTorrent Dataset and Tool Suite. 2013 10th Working Conference on Mining Software Repositories (MSR), 2013, pp. 233–236, doi: 10.1109/MSR.2013.6624034.
- [30] STEFAN, P.—HORKY, V.—BULEJ, L.—TUMA, P.: Unit Testing Performance in Java Projects: Are We There Yet? Proceedings of the 8th ACM/SPEC International Conference on Performance Engineering (ICPE '17), ACM, 2017, pp. 401–412, doi: 10.1145/3030207.3030226.
- [31] SULÍR, M.—BAČKOVÁ, M.—MADEJA, M.—CHODAREV, S.—JUHÁR, J.: Large-Scale Dataset of Local Java Software Build Results. Data, Vol. 5, 2020, No. 3, Art.No. 86, doi: 10.3390/data5030086.

- [32] PHAM, R.—SINGER, L.—LISKIN, O.—FILHO, F. F.—SCHNEIDER, K.: Creating a Shared Understanding of Testing Culture on a Social Coding Site. 2013 35th International Conference on Software Engineering (ICSE), 2013, pp. 112–121, doi: 10.1109/ICSE.2013.6606557.
- [33] ZHANG, Y.—LO, D.—KOCHHAR, P. S.—XIA, X.—LI, Q.—SUN, J.: Detecting Similar Repositories on GitHub. 2017 IEEE 24th International Conference on Software Analysis, Evolution and Reengineering (SANER), 2017, pp. 13–23, doi: 10.1109/SANER.2017.7884605.
- [34] JIANG, J.—LO, D.—HE, J.—XIA, X.—KOCHHAR, P. S.—ZHANG, L.: Why and How Developers Fork What from Whom in GitHub. *Empirical Software Engineering*, Vol. 22, 2017, No. 1, pp. 547–578, doi: 10.1007/s10664-016-9436-6.
- [35] KUHN, A.—VAN ROMPAEY, B.—HAENSENBERGER, L.—NIERSTRASZ, O.—DEMEYER, S.—GAELLI, M.—VAN LEEMPUT, K.: JExample: Exploiting Dependencies Between Tests to Improve Defect Localization. In: Abrahamsson, P., Baskerville, R., Conboy, K., Fitzgerald, B., Morgan, L., Wang, X. (Eds.): *Agile Processes in Software Engineering and Extreme Programming (XP 2008)*. Springer, Berlin, Heidelberg, Lecture Notes in Business Information Processing, Vol. 9, 2008, pp. 73–82, doi: 10.1007/978-3-540-68255-4_8.
- [36] Pearson's Correlation Coefficient. In: Kirch, W. (Eds.): *Encyclopedia of Public Health*. Springer, Dordrecht, 2008, p. 1090, doi: 10.1007/978-1-4020-5614-7_2569.



Matej MADEJA graduated (M.Sc.) at the Department of Computers and Informatics of the Faculty of Electrical Engineering and Informatics at the Technical University of Košice in 2017. He defended his Master's thesis in the field of informatics. Currently, he is Ph.D. student in the same department. His research is focused on the improvement of program comprehension efficiency, source code testing techniques, and teaching of programming.



Jaroslav PORUBÄN is Professor and the Head of Department of Computers and Informatics, Technical University of Košice, Slovakia. He received his M.Sc. degree in 2000 and his Ph.D. in 2004, both in computer science. Since 2003 he is a member of the Department of Computers and Informatics at the Technical University of Košice. Currently, the main subject of his research is computer language engineering focused on the design and implementation of domain-specific languages and computer language composition and evolution.



Michaela BAČKOVÁ is Assistant Professor and the Head of the Information Systems Laboratory at the Department of Computers and Informatics, Technical University of Košice, Slovakia. She received her Ph.D. in computer science in 2014. Since 2014 she is a member of the Department of Computers and Informatics at the Technical University of Košice. Currently, the main subject of her research is UX, HCI and usability while focusing on the domain-related terminology in user interfaces (domain usability). Her interest is also in software languages and innovations in the teaching process.



Matúš SULÍR is Assistant Professor at the Department of Computers and Informatics, Technical University of Košice, Slovakia. At the same university, he graduated with his Master's degree in computer science in 2014 and obtained his Ph.D. in 2018. His research is focused on program comprehension, particularly on the integration of run-time information with source code, attribute-oriented programming, and debugging. He is also interested in empirical studies in software engineering.



Jan JUHÁR is Researcher at the Department of Computers and Informatics, Technical University of Košice. He received his Ph.D. in computer science in 2018. Since 2018 he is a member of the Department of Computers and Informatics at the Technical University of Košice. His research focuses on program comprehension, programming tools, source code metadata, and program projections.



Sergej CHODAREV is Assistant Professor at the Department of Computers and Informatics, Technical University of Košice, Slovakia. He received his Master's degree in 2009 and his Ph.D. in 2012, both in computer science. The subject of his research includes domain-specific languages, metaprogramming, and software engineering.



Filip GURBÁL is Ph.D. student at the Department of Computers and Informatics, Technical University of Košice, Slovakia. He graduated at the university in computer science in 2020. He is a member of the Computer Network Laboratory at the Technical University of Košice. The subject of his research is improving program comprehension using methods and tools. He also focuses on software testing methods and tools.