# Are Unit and Integration Test Definitions Still Valid for Modern Java Projects? An Empirical Study on Open-Source Projects

Fabian Trautsch[a,*], Steffen Herbold[a], Jens Grabowski[a]

[a]*Institute of Computer Science, University of Goettingen, Germany*

## Abstract

**Context:** Unit and integration testing are popular testing techniques. However, while the software development context evolved over time, the definitions remained unchanged. There is no empirical evidence, if these commonly used definitions still fit to modern software development.

**Objective:** We analyze, if the existing standard definitions of unit and integration tests are still valid in modern software development contexts. Hence, we analyze if unit and integration tests detect different types of defects, as expected from the standard literature.

**Method:** We classify 38782 test cases into unit and integration tests according to the definition of the IEEE and use mutation testing to assess their defect detection capabilities. All integrated mutations are classified into five different defect types. Afterwards, we evaluate if there are any statistically significant differences in the results between unit and integration tests.

**Results:** We could not find any evidence that certain defect types are more effectively detected by either one of the test types. Our results suggest that the currently used definitions do not fit modern software development contexts.

**Conclusions:** This finding implies that we need to reconsider the definitions of unit and integration tests and suggest that the current property-based definitions may be exchanged with usage-based definitions.

---

[*]Corresponding author
*Email address:* `trautsch@cs.uni-goettingen.de` (Fabian Trautsch)

## 1. Introduction

Unit testing as well as integration testing are two popular testing techniques that are used in industry and academia. There exist a lot of scientific literature, which is focused on unit testing [1, 2, 3] or integration testing [4, 5, 6]. Furthermore, there exist several development models in which both practices are integrated, e.g., the V- Model [7] or the waterfall model [8].

The definitions for unit and integration tests, like the ones of the Institute of Electrical and Electronics Engineers (IEEE) [9] or International Software Testing Qualification Board (ISTQB) [10], are decades old. However, the software development context has evolved over time, while these definitions remained unchanged. For example, major contributions in the field of software testing that pushed this evolution were the development of xUnit frameworks (e.g., JUnit [11]), which are nowadays used to execute all types of tests including unit, integration, and system tests, or Continuous Integration (CI) systems (e.g., Travis CI [12], Jenkins [13]). Those frameworks and systems changed the way in which developers test their software. This is highlighted by a proposal that is spread within the development community [14]. The idea is to change a software testing paradigm: instead of testing a software thoroughly on unit level with a few integration tests, developers propose to test the software mostly on integration level with only a few unit tests. The reasons for this are manifold like, e.g., that integration tests are more realistic [14] or more effective [15, 16] than unit tests.

Hence, we can clearly see an evolution of the development context, but the definitions were not adapted over the time. Therefore, we hypothesize that the current definitions, that are used in current textbooks and software testing certifications (e.g., ISTQB certified tester [10]), do not fit to modern software development contexts. Our prior work [17] presents us indications that support

2

this hypothesis. We have shown that developers do not classify their tests according to the definitions of the IEEE or ISTQB. While this shows that the definitions are not used in practice, it does not gives us any indications regarding the merit of the current definitions. Since the main purpose of testing is to reveal defects, we focus on comparing the defect detection capabilities of unit and integration tests in this paper. Software testing textbooks state [18, 19, 20], that there should be a difference between unit and integration tests in terms of types of defects that they detect. Hence, we evaluate if unit or integration tests are more effective in detecting certain program defect types. If we do not find a difference, we would have another indication that the current definitions of unit and integration tests do not fit the modern software development context.

To steer our research, we defined one central Research Question (RQ):

- **RQ:** Are unit or integration tests more effective in detecting certain program defect types?

To answer our RQ, we have a look at the effectiveness of unit and integration tests in terms of their defect detection capabilities. For this, we first classify tests into unit and integration tests and afterwards collect data about their defect detection capabilities using mutation testing. Then, we perform an in-depth analysis by assessing which types of defects are found by the two different types of testing and if there are any differences in their effectiveness of detecting these defect types. The contributions of this paper can be summarized as follows:

- an in-depth analysis of the defect detection capabilities of unit and integration tests.

- an in-depth analysis of which defect types are found by unit and integration tests.

- a data set, which includes detailed information about the defect detection capabilities of 38782 test cases for 17 projects.

3

- a software test analysis framework, which can be used for further research and enable other researchers to contribute to the body of knowledge of empirical software testing research.

The remainder of this paper is structured as follows. In Section 2 we lay the foundations for this paper and discuss related work. Section 3 describes our research methodology, including the approach for our research question, our data collection and analysis procedures, our results, and the replication kit of this paper. We then discuss the results in Section 4. Section 5 presents the threats to validity to our study, including construct, internal and external threats. Finally, we conclude our paper and describe future work in Section 6.

## 2. Foundations and Related Work

In this section, we present the terminology that is used along this paper together with a summary of the literature that is related to our research.

### 2.1. Definitions

The IEEE standard ISO/IEC/IEEE 24765-2010 [9] defines the most important vocabulary for the software engineering world. In the following, we present the definitions that are important for this paper based on this standard.

**Definition 1 (Unit).** *1. a separately testable element specified in the design of a computer software component. 2. a logically separable part of a computer program. 3. a software component that is not subdivided into other components. [...]. [9]*

**Definition 2 (Unit Test).** *[...] 3. test of individual hardware or software units or groups of related units [9]*

**Definition 3 (Integration Test).** *1. the progressive linking and testing of programs or modules in order to ensure their proper functioning in the complete system. [9].*

4

## 2.2. Test Type Classification

There are several approaches present in the literature that try to classify tests into different classes. In the short paper by Orellana et al. [21], the classification is based on the usage of Maven plugins in project builds. If test classes are executed by the Maven SureFire plugin [22] they are classified as unit tests, but if they are executed by the Maven FailSafe plugin [23] they are classified as integration tests. Hence, they reuse the classification of the developers. This can be problematic, as our earlier work [17] highlights, that the classification into unit and integration tests, as done by the developers, is not always in line with the definitions. Furthermore, this kind of classification does not allow a fine-grained analysis on the method level. Another problem is, that not all projects differentiate their tests based on the executed Maven plugins. For example, none of the projects that we have used in our case study use the Maven FailSafe plugin, while some tests of the projects are indeed integration tests.

Another approach that classifies tests into different types was proposed by Kanstrèn [24]. Kanstrèn proposes a dynamic approach that uses aspect oriented programming to calculate the test granularity of each test, which refers "to the number of units of production code included in a test case..."[24]. To achieve this, his approach calculates the number of methods that are covered by each test. Afterwards, the results can be summarized to determine, if, e.g., a system was tested only on low-level (i.e., many tests that executed a small amount of methods) or only on high-level.

In our earlier work [17] we determined how many tests are unit tests, according to the definitions of the ISTQB and IEEE for 10 Python projects. We proposed a static analysis approach in which we assessed the number of imported modules in a test to find unit tests. Furthermore, we compared the intention that the developers had (i.e., if they wanted the test to be a unit test) with the outcome of the classification analysis.

We improved our test classification approach in several ways. We switched from a static analysis approach to a dynamic one to resolve the drawbacks of our former approach, i.e., we can now directly determine if a unit is used

(and not only imported) in a test case and the approach is robust against the usage of mocking frameworks inside a project. Our new approach is using test coverage data to classify tests into unit and integration tests, which was not possible before. Hence, our approach only needs the coverage data of a test execution run in contrast to the work of Orellana et al. [21]. Furthermore, our approach makes a clear separation between unit and integration tests based on the definition of the IEEE [9]. This is in contrast to the work of Kanstrèn [24], which only provides a general view on the granularity of the tests. Furthermore, our analysis is more fine-grained than before. Instead of classifying a whole test class/module we are able to classify each method separately.

*2.3. Assessment of Defect Detection Capabilities via Mutation Analysis*

The assessment of the defect detection capabilities of a set of test cases is done by checking how many defects the test suite can find. There exist several approaches in the literature that can be used. Recently, Oreallana et al. [21] assessed if unit tests expose more defects than integration tests by using the Travis Torrent dataset [25]. This data set contains information about project builds, including which test cases failed during the build. The authors of the paper used this information to determine the number defects that were exposed. The assumption in this kind of analysis is that every test that failed during the build of a project exposed a defect. The problem with this approach is that tests are often executed before committing changes to a Version Control System (VCS). Hence, all defects that might be found *before* committing the changes are not recorded. Moreover, this kind of technique is focused on defects that broke the build and therefore the number and type of assessed defects is very limited. First, because tests in a CI system could also fail due to wrong commit behavior (e.g., the developer forgot to commit a certain change to a class). Second, with this technique we could only assess defects that are detected by the CI system. Defects that might get fixed before committing to the CI system are missing. Hence, to answer our RQ we require a controlled environment, where we are able to exclude as many outside influences as possible that might

6

have an influence on the results. We achieve this through mutation analysis.

Mutation analysis provides a systematic way for introducing defects into a software [26]. In mutation analysis we integrate mutants (defects) into the program and check if these are found by test cases. If a test case detects a mutant, i.e. the test case fails, we call the mutant *killed*. Mutants that survive the execution of the test case are called *live* [26]. The *mutation score* is defined as a ratio of the number of killed mutants to the total number of mutants. Unfortunately, some of the generated mutants produce a program which behaves equivalent to the original one [27, 28] and can therefore not be killed. These mutants are called *equivalent*. However, the detection of such mutants is an instance of an undecidable problem [19]. Mutation testing can be very time consuming, as a large number of defect-seeded programs must be tested.

The use of mutation analysis for the assessment of the defect detection capabilities of tests, has the underlying assumption that the mutants can construct program failures that are similar to the ones that are created through real defects. There are several studies that provide indications that this is the case. One of the first studies that investigated the relationship between mutants and real defects was done by Daran et al. [29]. They found that injected mutants could produce failures and program data states that are similar to those produced by real defects. Similar to that, Andrews et al. [30, 31] concluded in their studies that the detection ratio of mutants are representative for defect detection ratios, i.e., there is a correlation between them. One of the most recent papers that highlight the relationship between mutants and real defects was done by Just et al. [32]. They found a strong correlation between mutant detection ratios and real defect detection ratios.

Nevertheless, there are some works that identified limitations to the conclusions presented above. Namin et al. [33] found that there is only a weak correlation between mutants and defect detection ratios. Chekam et al. [34] concluded that there is a strong connection between an increase of the mutation score and defect detection, but only at higher mutation score levels. A very recent paper by Papadakis et al. [35] highlight, that there is only a weak corre-

lation between the mutation score and defect detection, if the test suite size is controlled for.

As we can see on the above papers, the question if mutation testing is an appropriate tool to assess the defect detection capabilities of tests can not be definitely answered, as there exist indications for a positive as well as a negative answer to this question in the literature. While we also make use of mutation testing in our work, we do not create test suites, but reuse the ones that are provided by the projects. Hence, some limitations, e.g. that the test suite size must be controlled for, are not applicable to our research.

### 2.4. Defect Classification

There are numerous studies on defect classification. One of the main differences between these studies is the data on which the classification is based. While some taxonomies need specification or design documents [36], others only need source code [37], or defect reports [38, 39].

A commonly used cause-driven taxonomy was proposed by Chillarege et al. [40] and is called Orthogonal Defect Classification (ODC). In ODC defects are classified into eight types based on their description about the symptoms, semantics, and root causes. Offutt et al. [41] proposed a model for the characterization of defects based on the syntactic and semantic size. Hayes et al. [36] presented a requirements-based defect analysis methodology, which the author also applied to NASA projects. Later, Hayes et al. [42] developed two taxonomies, where the first classifies code modules (e.g., into data-centric, controller, view, ...) and the second code defects (e.g., into data, interface, computation, ...). Xia et al. [38] used the descriptions available in defect reports to classify defects into two defect trigger categories (Bohrbug and Mandelbug) via natural language processing techniques. Tan et al. [39] created a taxonomy in which defects are classified based on three dimensions: root cause (of the defect), impact (failure caused by the defect), and component (location of the defect).

The main problem with the approaches described above is that the data that is needed to classify the defects are often not available, e.g. specifications

8

or design documents. This is especially true for open source projects, as these projects follow a special development process [43]. Additionally, the creation of a link between the classified defect and its representation in the source code is often hard to achieve.

Recently, Zhao et al. [37] presented an approach that can overcome the problems described above. They classify defects based on the change that was made to fix the defect. For this, they adapted the classification of [42] and created a tool for the C language that needs two versions of the source code of the program (defective and clean). Afterwards, the differences (changes) between these two versions are extracted and change patterns are detected. Based on these change patterns a change type classification is created, which is also the classification of the defect. They created five different categories with overall nine subcategories. The **data** category comprises of Changes on Data Declaration/Definition (CDDI) statements. This type of change indicate that a data-related defect occurred, e.g., if the declared type of a variable is changed from *int* to *float*. The **computation** category includes Changes on Assignment Statements (CAS). This includes the addition or deletion of assignment statements, or the modification of equations. Computation-related defects are defects that can lead to a wrong assignment of a variable, which are then often fixed by CAS. **Interface**-related defects are caused "by wrong definition or faulty function dependency on other functions." [37]. For example, if a function is called with an incorrect amount of parameters. Defects that are related to interface issues are fixed by Changes on Function Declaration/Definition (CFDD) or Changes on Function Call (CFC). The **logic/control** category comprises of Changes on Loop Statements (CLS), Changes on Branch Statements (CBS), and Changes on Return/Goto Statements (CRGS). Hence, defects that occur in these statements, e.g., changing a $<$ to a $>=$ in an if statement, "may cause the incorrect execution sequence or an abnormal state" [37]. Therefore, a defect-fixing change in these statements signals that a logic/control-related defect occurred. Every other change that could not be classified into the categories above is then subsumed under the **Other** category. For example, Changes on

Preprocessor Directives (CPD), which is a change that can occur in programs using the C language.

Our approach reuses the classification scheme of Zhao et al. [37] as it is tightly connected to the source code. This connection to the source code is needed to classify our generated mutants, as these mutants are not regular defects that exhibit, e.g., an issue report. While we reuse the logic behind every sub category our tool only outputs the main category in which a classified defect resides in (e.g., "Interface"). Furthermore, in contrast to [37], our tool is working for Java projects.

## 3. Research Methodology

Within our RQ, we want to assess if unit or integration tests are more effective in detecting certain program defect types. Our expectation would be that integration tests detect more interface problems, while unit tests are more focused on finding computation defects or problems with the control flow of a program. This is also stated in software testing textbooks [18, 19, 20]. Hence, we evaluate if there are indeed defect types that are only (or mostly) detected by a certain test type and assess how effective these test types can detect certain defect types. These results could help to assess if the current definitions of the IEEE still fit to modern software development contexts.

Overall, our approach is as follows. We classify all test cases into unit and integration tests. Additionally, we collect data on the defect detection capabilities of each test case by using mutation analysis. We generate many different defective versions of the code and check which tests are able to detect the integrated defects. Afterwards, we classify our introduced defects based on the changes that were made to fix the defect using the taxonomy from Zhao et al. [37]. Finally, we evaluate which defect types are found by which test type. Within this section, we describe our case study design, including the subject selection, data collection, and analysis procedures.

*3.1. Subject Selection*

We defined the following inclusion criteria to select our study subjects.

- **Projects must be a library or a framework.** Libraries and frameworks should be included or used by other programs and are not executed by themselves. Therefore, tests that execute the whole system (system tests) are less likely to occur. Hence, by focusing on libraries and frameworks we are reducing the risk of misclassifying a system test as an integration test, as our tooling infrastructure is currently not capable to differentiate between these test types.

- **Projects must have a minimum of 1000 commits and are at least 2 years old.** This ensures, that we only include mature projects.

- **Projects must use Java as programming language, Maven [44] as build system, and JUnit [11] or TestNG [45] as test driver.** This is a limitation of our current tooling infrastructure.

Moreover, one exclusion criterion is defined.

- **Projects should not be focused on Android alone.** This work concentrates on pure Java projects. While Android projects also use Java as programming language, there are several differences in testing those projects in contrast to pure Java projects.

We applied these criteria on two different data sources. First, the list of Borges et al. [46], which classified the most popular 5000 GitHub repositories (language-independent) into six different categories (application software, system software, web libraries and frameworks, non-web libraries and frameworks, software tools, and documentation). Second, we used a list of the most popular projects created by the maven repository [47][1]. Overall, 41 projects fit to our criteria from which we randomly selected 17 projects for our study.

---

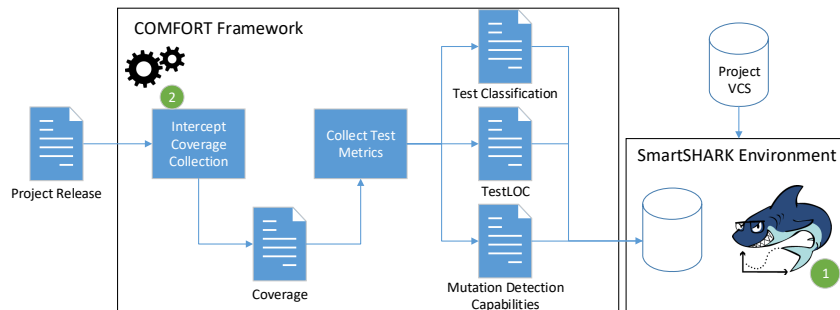[1]As available on the time we performed our case study (January 2018).

**Figure 1** Logical overview of our data collection process. Step 1 is the collection of meta- data using the SmartSHARK environment [48]; Step 2 is the automatic collection of test metrics, including the mutation detection capabilities of all tests.

Table 1 gives an overview of the chosen projects together with their characteristics. It shows the project name, the release of the project that was used in the case study, the overall number of commits (till the stated release), the number of Java classes for the release, and the number of test cases that are executed by the build script of the release. All projects can be found on GitHub. The number of classes and test cases in this table highlight that the projects that we included in our case study are non-trivial projects.

From the sampled projects we included the latest minor release, if available. If this release was not available, we took the next possible release (e.g., release 1.11.1 of *jsoup*, as the release 1.11.0 is not available). For two projects (i.e., *commons-beanutils* and *fastjson*) the latest minor release was longer than four years ago. Hence, we decided to use the most recent release to prevent compatibility issues with our infrastructure.

*3.2. Data Collection Procedures*

Figure 1 gives a logical overview of our data collection process for our case study. As a first step, we are collecting meta-data about each project using the SmartSHARK environment [48]. Hence, we collect data from the VCSs of the projects, as these are later on used in our data collection process to integrate the different collected data. Afterwards, in the second step, we calculate different

| Project | Release | #Commits | #Java Classes | #Test Cases |
|---|---|---|---|---|
| commons-beanutils | 1.9.3 | 1128 | 370 | 1197 |
| commons-codec | 1.11 | 1677 | 166 | 874 |
| commons-collections | 4.1 | 2852 | 903 | 6637 |
| commons-io | 2.5 | 1868 | 308 | 1150 |
| commons-lang | 3.7 | 5106 | 711 | 4064 |
| commons-math | 3.6 | 5819 | 2596 | 6488 |
| druid | 1.1.0 | 5178 | 3785 | 4132 |
| fastjson | 1.2.41 | 2579 | 5214 | 4176 |
| google-gson | 2.8.0 | 1306 | 719 | 1014 |
| guice | 4.1 | 1513 | 2078 | 701 |
| HikariCP | 2.7.0 | 2493 | 154 | 120 |
| jackson-core | 2.9.0 | 1323 | 268 | 774 |
| jfreechart | 1.5.0 | 3622 | 1041 | 2175 |
| joda-time | 2.9 | 1913 | 527 | 4176 |
| jsoup | 1.11.1 | 1106 | 304 | 593 |
| mybatis-3 | 3.4.0 | 1816 | 1205 | 1053 |
| zxing | 3.3.0 | 3282 | 408 | 401 |

**Table 1** Selected projects with their characteristics.

test metrics for each analyzed test case. This step is automated by our Collection of Metrics for Tests (COMFORT) framework. In a first step, we execute the tests of a project release and intercept the coverage collection to collect the coverage on a *per test case level*. The resulting coverage is afterwards used to calculate different test metrics for each test case: the test classification into unit and integration test (see Section 3.2.1), the Test Lines of Code (TestLOC) (see Section 3.2.2), and the mutation detection capabilities (see Section 3.2.3). In the following, we describe each approach that we followed to calculate these metrics and the concrete implementation that we designed. Furthermore, we describe our approach for the collection of the defect type in Section 3.2.4.

*3.2.1. Assigning the Test Classification*

Within this section, we explain our approach to assign a test class (i.e., unit or integration test) to the tests of the projects. We base our our separation on the IEEE definitions presented in Section 2.1, because the ISTQB definition (i.e., a unit test tests only *one* unit) is a subset of it. For example, all ISTQB unit tests are also IEEE unit tests by definition. Furthermore, our prior work [17] shows that the ISTQB definitions, are too restrictive for modern software systems, where *logical* software units often consist of several classes.

According to the definitions of the IEEE (see Section 2.1), a *test* in Java is a class, which consists of several test methods. These test methods are *test cases*, which are e.g., annotated with the "@Test"-JUnit annotation [11]. A unit is a *class* in Java, as it is a logically separable part of a Java program, which is not further subdivided into other components. This is in line with the literature [49]. A unit test is a test case that tests "...software units or groups of related units" [9]. In Java, related units are put into one package, as the official Java documentation states: "A package is a namespace that organizes a set of related classes and interfaces." [50]. Hence, if we apply the IEEE definition to Java a unit test is a test that tests only units from within one package (i.e., related units). On the other hand, an integration test tests units from more than one package. Hence, we do not classify the *intent* with which a test is created, but the actual type of the test according to the IEEE definitions.

To foster the understanding of the IEEE definitions, we selected two real world examples from the *commons-io* project [51]. Listing 1 shows an IEEE unit test and an IEEE integration test. The upper part of this listing presents a typical unit test, where the correct workings of a function is tested. Here, the *getPrefix* function of the *FilenameUtils* class is tested. More precisely, this test checks if null bytes are part of the string. As this tests only calls one unit (i.e., the *FilenameUtils* class) and the *FilenameUtils* class does not call other units, this test gets classified as a unit test. On the other hand, the bottom part of the listing shows an IEEE integration test. In this case, the test checks if

14

the *ByteArrayOutputStream* class of *commons-io* works as intended, if it is used within the *copy* function of the *CopyUtils* class. As the *ByteArrayOutputStream* is from the *org.apache.commons.io.output* package and the *CopyUtils* class from the *org.apache.commons.io* package, this tests gets classified as integration test.

```java
1   // IEEE Unit Test
2   [...]
3   package org.apache.commons.io;
4   [...]
5
6   @Test
7   public void testGetPrefix_with_nullbyte() {
8     try {
9       assertEquals("~user\\", FilenameUtils.getPrefix("~u\
            u0000ser\\a\\b\\c.txt"));
10    } catch (IllegalArgumentException ignore) {
11    }
12  }
13
14  /****************************************************/
15
16  // IEEE Integration Test
17  [...]
18  package org.apache.commons.io;
19  [...]
20  import org.apache.commons.io.output.ByteArrayOutputStream;
21  [...]
22
23  @Test
24  public void copy_byteArrayToOutputStream() throws Exception {
25    final ByteArrayOutputStream baout = new ByteArrayOutputStream
            ();
26    final OutputStream out = new YellOnFlushAndCloseOutputStream(
            baout, false, true);
```

15

```
27

28    CopyUtils.copy(inData, out);

29

30    assertEquals("Sizes differ", inData.length, baout.size());
31    assertTrue("Content differs", Arrays.equals(inData, baout.
          toByteArray()));
32  }
```

**Listing 1** Examples for an unit test (top) and an integration test (bottom) from the *commons-io* project [51].

On the implementation side, we use the recorded per test case coverage data for the test class assignment. A unit counts as covered, if at least one method of it was executed by the test. Furthermore, we only record the coverage of units that are *within* the project, i.e., calls to other libraries or the Java standard Application Programming Interface (API) or mocking frameworks are filtered out. In addition, we filter out every test class that might be inside the coverage data, as we want to base our test class assignment only on covered production classes. Therefore, we check the file path for each covered class and if the path has the term "test"[2] in it, we filter it out from the data. This has another advantage: self-made mocks that might exist in the project are filtered out from the coverage data. Within this filtered coverage data, we check how many "related units" (i.e., different packages) are covered by a test case. If only one logical unit is covered, then the test case is classified as unit test. Otherwise, it is classified as an integration test.

*3.2.2. Collecting the TestLOC*

We need to acquire the number of TestLOC for our RQ to normalize the results. This number is calculated based on the coverage data that is also used for the test classification. We sum up the number of covered lines for each

---

[2]In Maven projects (like we use in our study) all tests and test related data is in the "src/test" folder.

test case, while we differentiate between covered production lines and covered test lines. This differentiation is done based on the path in which the covered class resides. If the path starts with "src/test" we know that the covered class is a test class, because it resides in the test folder, and is a production class otherwise.

In contrast to just counting Lines of Code (LOC) or Logical Lines of Code (LLOC) this approach has the advantage that we consider all executed test code than just the lines of the test case itself. Otherwise our results could be biased. For example, tests that have longer set up methods (e.g., to bring several units into a testable state) but less LOC would show better results after the normalization than tests where the whole set up is integrated into the test case itself.

*3.2.3. Collecting the Mutation Detection Capabilities*

We collect the mutation data for every test case that was executed (i.e., is available in the coverage data) by using PIT [52] within our framework. We execute PIT for every test case separately, i.e., for each test case that was covered. Each test case is run against all generated mutations. Afterwards, the results of PIT are parsed. If a test case failed when it is executed alone, e.g., because this test case depends on the execution of other test cases, we excluded it and no mutation data is collected. As the collection of the mutation data is very computation intensive (i.e., our case study experiment consumed about 35 months of CPU time) we executed it on the High Performance Computing (HPC) cluster hosted by the datacenter of our university[3].

We followed the best practices on using mutation testing in controlled experiments by Papadakis et al. [26]. Hence, we also report on each decision regarding the mutation analysis, as advised by Papadakis et al. [26] in the following:

- **Mutation Testing Tool:** We used the paper by Kintis et al. [53] as a basis for our decision regarding the mutation testing tool. We decided

---

[3]See: http://gwdg.de.

to use PIT [54] for our analysis due to several reasons. PIT possess a good defect-revelation ability [53], which is an indicator of its suitability for our experiments. Only $PIT_{RV}$ [52, 55] was able to reveal more defects (115:122). Nevertheless, we decided against $PIT_{RV}$, because the number of generated *equivalent* mutants is substantially higher for $PIT_{RV}$ in contrast to PIT [53]. In fact in the experiment by Kintis et al. [53], $PIT_{RV}$ generated about 88.74% more equivalent mutants than PIT. As equivalent mutants are a substantial threat to the validity of our study, we decided to use PIT instead of $PIT_{RV}$. We even contributed to the development of PIT by creating a pull request that got accepted and integrated into the version 1.3.2 of PIT. This addition was required for our work and allows PIT users to filter the test cases of tests against which the mutations are challenged.

- **Mutant Redundancy:** Mutant redundancy might have a large impact on the validity of our study. Hence, it is important to care for this kind of thread. We generated a lot of mutants for our case study ($> 500.000$). Checking them all by hand to detect duplicate mutants is a task that would take a substantial amount of time and is error prone. To mitigate this threat, we create the set of disjoint mutants, which only includes mutants that are not killed collaterally by most of the test cases [53, 56].

- **Mutant Selection:** We decided to use all mutation operators that are provided by PIT [54], including operators that are also used in integration mutation testing approaches [57, 58]. The reason for this is that we wanted to generate a huge set of mutants so that we reduce the possibility that we only generate mutants of low quality or mutants that favor only one specific type of test. There is also empirical evidence that supports this approach [53].

- **Test Suite Choice and Test Suite Size:** The test suite choice and its size is set, as we reuse the test suites that were created by the developers of the projects.

- **Clean Program Assumption:** The general problem that the Clean Program Assumption (CPA) describes is that test suites are assessed on the mutated program instead the original one for which they were created [34]. While we know of the CPA, we do not need to take it into account for our experiment, as we do not compare different testing techniques or rely on the coverage measured on the clean program.

- **Multiple Experimental Repetitions:** The calculations for our case study are only done once, as we do not have a component in it that make stochastic choices. The mutants that are generated for each test case and the result (i.e., mutant was killed or not) are not affected by any stochastic process. An exception is the creation of the disjoint mutant set, as the used algorithm is not deterministic. Hence, we repeated the generation of the disjoint mutant set 10 times.

- **Presentation of the Results:** As described above, our results are presented on the test case level.

*3.2.4. Collecting the Defect Classification*

To gather the defect classification for each integrated mutation, we created the BugFixClassifier [59]. The BugFixClassifier is able to extract changes between two files using CHANGEDISTILLER by Fluri et al. [60] and classify the detected changes based on the approach by Zhao et al. [37]. CHANGEDISTILLER [60] is a tool that extract changes between two source code files based on the comparison of the Abstract Syntax Trees (ASTs) of both files. For each change that CHANGEDISTILLER detects, it outputs the change type (as defined by Fluri et al. [60]), which entity was changed (e.g., an if-statement or a method), and the parent entity of the changed entity (e.g., the initialization part of a for loop). We decided to reuse CHANGEDISTILLER, as it provides us with fine-grained source code changes that can be mapped onto the defect types as defined by Zhao et al. [37]. Hence, we do not only store each generated mutation, but also the type of defect that was integrated.

Zhao et al. [37] determines the class of a defect based on the change that was made to fix the corresponding defect. The only change that we made to the classification schema is that we excluded the category CPD, because pre-processor directives are not available in Java. We reuse the detailed change types of Fluri et al. [60] and provide a mapping between them and the defect classes defined by Zhao et al. [37]. This mapping is only provided for the upper categories of our classification scheme (i.e., data, computation, interface, logic/control, others), as this level of detail is sufficient for our purpose. Table A.7 in Appendix A shows the mapping between the change types of Fluri et al. [60] and the defect classes of Zhao et al. [37]. Unfortunately, some change types of Fluri et al. [60] cannot be directly mapped onto a defect class. For these change types we need to consider the type of the changed entity and/or the type of the parent entity (e.g., METHOD, FOR_INIT). Table A.8 in Appendix A shows the conditions that need to be fulfilled for a change together with the defect class that is then assigned.

In general, the original non-defective and the defective source code is needed for the defect classification. The changes between both are then extracted and classified based on the rules above. Unfortunately, our used mutation testing framework does not offer the possibility to get the mutated source code, as the mutation is done on the byte-code of the original code [52]. Nevertheless, some of the applied mutations can be directly mapped to a defect class. This mapping is presented in Table 2. For mutation operators that can not be directly mapped we used the following approach. The mutation testing framework outputs the mutation operator that is used and the line in which it was used. Hence, we integrate the corresponding defect (based on the mutation operator) into the original source code in the specified line. Afterwards, we extract the changes between the original and the now defective source code to get the class of the defect that was integrated by the mutation operator.

## 3.3. Data Analysis and Results

In the following sections, we explain our analyzed data sets, our analysis procedures, as well as our results for our RQ.

### 3.3.1. Data Sets

Table 3 shows the number of unit and integration tests for each project release that we analyzed, together with the number of not classifiable tests, the sum of the TestLOC for unit and integration tests, the number of unique mutants that are generated, and the number of analyzed tests. Overall, this table highlights that nearly all projects have more integration than unit tests. Nevertheless, for some projects our approach was not able to classify all test cases. We looked into all not classifiable test cases and found several reasons for this, e.g., tests that are empty, tests that directly return without executing any production code, or test cases that only test constants of classes. Furthermore, Table 3 highlights the number of generated unique mutants[4] together with the number of analyzed tests. The number of analyzed tests can be lower than the number of all test cases of the projects, as our mutation testing tool might not be able to run the test alone (e.g., test cases that fail if they are executed alone).

---

[4]Not all mutants are generated for each test case as our mutation testing tool pre-selects mutations against which the test case should run by using the coverage data of the test case [52].

| Mutation Operator | Defect Class |
|---|---|
| **AP**: *Argument Propagation* | Interface |
| **BFR**: *Boolean False Return* | Logic/Control |
| **BTR**: *Boolean True Return* | Logic/Control |
| **CB**: *Conditionals Boundary* | - |
| **CC**: *Constructor Calls* | Interface |
| **EOR**: *Empty Object Return* | Logic/Control |
| **I**: *Increments* | - |
| **IC**: *Inline Constant* | Data |
| **IN**: *Invert Negatives* | - |
| **M**: *Math* | - |
| **MV**: *Member Variable* | Computation |
| **NR**: *Naked Receiver* | Interface |
| **NC**: *Negate Conditionals* | - |
| **NVMC**: *Non Void Method Calls* | Interface |
| **NR**: *Null Return* | Logic/Control |
| **PR**: *Primitive Return* | Logic/Control |
| **RC**: *Remove Conditionals* | Logic/Control |
| **RI**: *Remove Increments* | - |
| **RS**: *Remove Switch* | Logic/Control |
| **RV**: *Return Values* | Logic/Control |
| **S**: *Switch* | Logic/Control |
| **VMC**: *Void Method Calls* | Interface |

**Table 2** Mapping between the used mutation operators and the defect class.

| Project | #UT | #IT | #NC | TestLOC (UT) | TestLOC (IT) | #Unique Mutants | #Analyzed Tests |
|---|---|---|---|---|---|---|---|
| commons-beanutils | 285 | 890 | 22 | 10966 | 44510 | 11310 | 1175 |
| commons-codec | 707 | 146 | 21 | 6160 | 1136 | 9059 | 853 |
| commons-collections | 1925 | 4005 | 707 | 63048 | 260958 | 26464 | 5930 |
| commons-io | 634 | 504 | 12 | 17345 | 13185 | 9593 | 1138 |
| commons-lang | 3507 | 471 | 86 | 36718 | 10841 | 36549 | 3978 |
| commons-math | 775 | 5709 | 4 | 10812 | 157902 | 113469 | 6484 |
| druid | 94 | 4037 | 1 | 585 | 71456 | 123835 | 4127 |
| fastjson | 315 | 3842 | 19 | 1968 | 41312 | 48289 | 4147 |
| google-gson | 215 | 797 | 2 | 2185 | 11115 | 8816 | 1012 |
| guice | 28 | 673 | 0 | 265 | 13526 | 10851 | 688 |
| HikariCP | 21 | 97 | 2 | 291 | 5848 | 4967 | 117 |
| jackson-core | 47 | 727 | 0 | 490 | 17567 | 34361 | 774 |
| jfreechart | 228 | 1946 | 1 | 2642 | 31240 | 96104 | 2174 |
| joda-time | 179 | 3975 | 22 | 2810 | 83739 | 32951 | 4153 |
| jsoup | 64 | 525 | 4 | 705 | 11772 | 14171 | 588 |
| mybatis-3 | 224 | 824 | 5 | 1449 | 28577 | 17126 | 1043 |
| zxing | 108 | 293 | 0 | 1564 | 13809 | 29592 | 401 |
| Overall | 9356 | 29461 | 908 | 160003 | 818493 | 627507 | 38782 |

**Table 3** Projects with their collected data, including Unit Tests (UT), Integration Tests (IT), and not classifiable tests (NC).

We create two different data sets for the analysis of our research question based on the collected data. These data sets represent different perspectives on the questions at hand.

- **ALL:** This data set consists of the test results for all generated mutations. With this data set we want to assess the defect detection capabilities of unit and integration tests for a large data set with many different defects that are integrated.

- **DISJ:** This data set consists of the test results for the set of disjoint mutants (see Section 3.2.3). With the analysis of this data set we can gain insights into the defect detection capabilities of unit and integration tests for defects that are hard to kill [56].

### 3.3.2. Analysis Procedure

For all of the in Section 3.3.1 presented data sets we executed the following analysis process.

1. **Calculate the number of detected defects:** We gather the number of all detected defects by summing up the number of defects that are detected by a unit or an integration test. We consider a mutation that is *killed* as a detected defect. The results for test cases that are executed with several different parameters, are combined and analyzed as one test case. As the algorithm applied to create the set of disjoint mutants is not deterministic, we repeat the analysis process using the disjoint mutant set 10 times and take the average of all 10 runs for the number of detected defects.

2. **Build the sum for each test type:** We divide the detected defects by their type and sum them up for unit and integration tests separately. This way, we can assess if unit or integration tests are more effective in detecting a certain kind of defect. We excluded the defect type *Other* from our analysis, as it does not represent a real defect type, but more a type of change that can not be classified as one of the other types (see

Section 2.4). Note, that we have prespecified the subgroup analysis for this paper. Hence, it is not a post hoc analysis, which would threat the validity of our study [61].

3. **Normalize by the number of TestLOC:** The resulting sums from the previous step are normalized by the number of TestKLOC to create scores. Hence, the result is the number of detected defects per 1000 TestLOC. We perform this normalization step because we want to include the effort that is done to create tests into our analysis. Hence, this normalization step is needed, as we would otherwise just compare the number of detected defects, which would ignore the effort (i.e., the TestLOC).

4. **Statistical Testing:** As a last step, we take the normalized values of each project and perform statistical tests to check if the differences between the unit and integration test scores are significant. We first check if our populations follows a normal distribution (applying the Shapiro-Wilk test [62]) and have equal variances (applying the Levene test [63]) to choose an appropriate significance test. Based on the results, we either perform a students t-test [64] or a Mann-Whitney-U test [65]. We use a significance level of $\alpha = 0.05$ for all of our statistical tests. As we are performing multiple statistical tests, which could increase the overall change of false discoveries (Type 1 Errors) [66], we need to apply corrections for multiple comparisons. We decided to use the Bonferroni correction [66] for all of our statistical hypothesis tests that we made. Overall, we use 8 statistical hypothesis tests: we use our **ALL** and **DISJ** data sets and check for differences in the scores for each defect type (i.e., COMPUTATION, DATA, INTERFACE, LOGIC/CONTROL), which results in eight different tests (two data sets * 4 defect types). Therefore our adjusted significance level is $\alpha^* = 0.05/8 = 0.00625$.

### 3.3.3. Results

Tables 4, and 5 show the scores (i.e., number of detected defects per TestKLOC) of unit and integration tests, separated by the type of defect that they

have detected. Additionally, the mean and standard deviation is shown for each column. Additional tables including the number of defects found by unit tests, integration tests, and both for each defect type are available in the supplementary material (see Section 3.4). Furthermore, the supplementary material includes Venn-Diagrams for each project to visualize these numbers.

Table 4 depicts the results for the data set **ALL**. It shows that if we separate the integrated defects by type, the mean scores of integration tests are higher than the mean scores of unit tests for each defect type except DATA defects, while they have a higher standard deviation. Hence, on average it seems that integration tests are more effective (i.e., scores are higher) for any defect type, except DATA defects. Furthermore, Table 4 also highlights the differences between projects. For some projects the results are as expected (i.e., integration tests detect interface defects more effectively, while other defects are more effectively detected by unit tests), while for other projects it is vice versa (e.g., *jackson-core*).

| | COMP. | | DATA | | INT. | | L/C | |
|---|---|---|---|---|---|---|---|---|
| **Project** | UT | IT | UT | IT | UT | IT | UT | IT |
| commons-beanutils | 3.83 | 2.29 | 6.29 | 3.30 | 38.39 | 32.44 | 69.31 | 44.69 |
| commons-codec | 71.27 | 22.01 | 134.42 | 66.02 | 252.44 | 344.19 | 389.12 | 448.94 |
| commons-collections | 2.24 | 1.25 | 2.98 | 1.11 | 8.31 | 7.05 | 20.08 | 13.62 |
| commons-io | 11.01 | 15.09 | 20.12 | 19.34 | 57.94 | 64.77 | 101.64 | 104.29 |
| commons-lang | 24.59 | 17.06 | 61.36 | 44.46 | 125.25 | 162.62 | 354.29 | 276.36 |
| commons-math | 28.39 | 38.68 | 41.90 | 76.47 | 86.57 | 129.81 | 195.15 | 215.08 |
| druid | 133.33 | 57.56 | 169.23 | 55.92 | 194.87 | 396.17 | 558.97 | 501.95 |
| fastjson | 31.50 | 50.47 | 66.57 | 109.65 | 63.01 | 213.72 | 263.72 | 328.86 |
| google-gson | 66.36 | 17.09 | 59.04 | 14.22 | 65.45 | 115.61 | 174.37 | 157.17 |
| guice | 11.32 | 37.93 | 18.87 | 31.05 | 226.42 | 223.94 | 362.26 | 264.82 |
| HikariCP | 37.80 | 30.61 | 65.29 | 28.73 | 151.20 | 125.68 | 109.97 | 163.13 |
| jackson-core | 87.76 | 153.07 | 136.73 | 154.32 | 293.88 | 195.42 | 561.22 | 589.97 |
| jfreechart | 36.34 | 76.47 | 71.16 | 92.93 | 75.32 | 290.01 | 281.98 | 593.41 |
| joda-time | 7.47 | 13.76 | 19.93 | 21.91 | 65.12 | 100.98 | 87.19 | 158.89 |
| jsoup | 5.67 | 32.45 | 19.86 | 58.78 | 41.13 | 280.33 | 114.89 | 372.07 |
| mybatis-3 | 15.87 | 13.89 | 18.63 | 12.25 | 193.24 | 119.15 | 204.28 | 116.21 |
| zxing | 41.56 | 81.83 | 190.54 | 211.46 | 203.32 | 282.42 | 435.42 | 611.12 |
| Mean | 36.25 | 38.91 | 64.88 | 58.94 | 125.99 | 181.43 | 251.99 | 291.80 |
| StDev | 35.45 | 37.67 | 58.56 | 56.96 | 86.22 | 110.24 | 168.34 | 197.19 |

**Table 4** Scores for unit and integration tests for the **ALL** data set, separated by defect type.

Table 5 highlights the results for the **DISJ** data set. It shows a different picture than Table 4. For the disjoint mutant set our results show that (on average) integration tests are less effective in detecting any type of defect than unit tests, i.e., the mean of integration test scores are lower than the mean of unit test scores for any of the defect types. But, the standard deviation for integration test scores are lower than the standard deviation of unit test scores for any of the defect types. Nevertheless, Table 5 also highlights, that some mutations are only detected by integration tests (e.g., interface defects for the

projects *fastjson*, *google-gson*, *guice*, *HikariCP*, *jackson-core*) or only unit tests (e.g., computation defects for the projects *commons-beanutils*, *commons-lang*).

| Project | COMP. UT | COMP. IT | DATA UT | DATA IT | INT. UT | INT. IT | L/C UT | L/C IT |
|---|---|---|---|---|---|---|---|---|
| commons-beanutils | 0.09 | 0.00 | 0.00 | 0.02 | 0.18 | 0.09 | 0.64 | 0.38 |
| commons-codec | 1.30 | 0.88 | 0.65 | 0.00 | 2.44 | 2.64 | 7.31 | 2.64 |
| commons-collections | 0.02 | 0.00 | 0.02 | 0.01 | 0.13 | 0.04 | 0.40 | 0.11 |
| commons-io | 0.35 | 0.15 | 0.23 | 0.08 | 1.56 | 0.15 | 1.50 | 0.76 |
| commons-lang | 0.44 | 0.00 | 0.93 | 0.18 | 2.31 | 0.46 | 8.14 | 0.46 |
| commons-math | 0.28 | 0.28 | 0.00 | 0.16 | 0.37 | 0.16 | 1.85 | 1.04 |
| druid | 0.00 | 0.64 | 6.84 | 0.34 | 6.84 | 1.75 | 17.09 | 2.70 |
| fastjson | 0.51 | 1.02 | 0.00 | 1.19 | 0.00 | 2.52 | 1.52 | 6.00 |
| google-gson | 0.92 | 0.09 | 0.00 | 0.27 | 0.00 | 0.09 | 0.92 | 0.81 |
| guice | 3.77 | 0.30 | 0.00 | 0.44 | 0.00 | 1.11 | 0.00 | 1.85 |
| HikariCP | 0.00 | 0.00 | 0.00 | 0.17 | 0.00 | 0.17 | 0.00 | 0.17 |
| jackson-core | 0.00 | 0.23 | 2.04 | 0.06 | 0.00 | 0.17 | 10.20 | 0.91 |
| jfreechart | 1.89 | 0.29 | 0.38 | 0.10 | 0.76 | 0.54 | 5.30 | 1.28 |
| joda-time | 0.00 | 0.12 | 0.00 | 0.16 | 0.71 | 0.57 | 1.42 | 2.67 |
| jsoup | 0.00 | 0.00 | 0.00 | 0.08 | 0.00 | 0.08 | 0.00 | 0.08 |
| mybatis-3 | 0.69 | 0.14 | 0.69 | 0.00 | 1.38 | 0.07 | 0.00 | 0.03 |
| zxing | 0.64 | 0.00 | 0.00 | 0.14 | 0.00 | 0.07 | 0.64 | 0.29 |
| Mean | 0.64 | 0.24 | 0.69 | 0.20 | 0.98 | 0.63 | 3.35 | 1.31 |
| StDev | 0.97 | 0.31 | 1.67 | 0.28 | 1.72 | 0.86 | 4.78 | 1.52 |

**Table 5** Scores for unit and integration tests for the **DISJ** data set, separated by defect type.

Table 6 presents the p-values of the significance tests that were performed between the scores of unit and integration tests for the data sets **ALL** and **DISJ** and each defect type. It highlights, that there are no statistical significant differences in the effectiveness of unit and integration tests for any defect type.

| Defect Type | ALL | DISJ |
|---|---|---|
| COMPUTATION | $p = .352$ | $p = .152$ |
| DATA | $p = .766$ | $p = .220$ |
| INTERFACE | $p = .112$ | $p = .278$ |
| LOGIC/CONTROL | $p = .531$ | $p = .267$ |

**Table 6** P-values of the significance tests that were performed between the scores of unit and integration tests for the data sets **ALL** and **DISJ** and each defect type.

> **Answer to our RQ:** Tables 4 and 5 highlight, that there are differences in the effectiveness between unit and integration tests for the different defect types if we compare the means in the tables. The results which test type is more effective for which defect type is not consistent over our data sets, indicating that unit tests are more effective in detecting "hard to kill" defects. Furthermore, the differences in the effectiveness are not significant, as Table 6 highlights. Our results also vary from project to project: there are some projects which seem to have more effective unit tests for certain defect types, while in others integration tests are more effective for the same defect type.

*3.4. Replication Kit*

To facilitate further insights and the replication of our study, we provide a replication kit [67]. The replication kit contains the following data:

- all source code used for the data collection, including the used versions of the COMFORT framework, the BugFixClassifier, and the used tools of the SmartSHARK environment;

- all source code used for the analysis of the collected data;

- a MongoDB database dump with all raw results, except the developer information; and

- additional visualizations of the results.

## 4. Discussion

The results of our research were rather unexpected for us, as it does not represents what we have learned and what we teach at our university. Overall, we could not find any significant differences in the effectiveness between unit and integration tests for any of the data sets that we tested. These results are interesting out of several perspectives.

**Education:** Academia, as well as organizations like the ISTQB, teach that integration and unit tests are equally important and find different types of defects (i.e., integration tests detect interface defects, whereas unit tests detect other kind of defects). They also highlight that a separation between unit and integration tests make sense and should be done. Nevertheless, our results show that there is no significant difference in the effectiveness between unit and integration tests for any of the defect types. Hence, it seems that unit and integration tests are even equally effective in detecting interface defects. In addition, we have seen that there are defects that are detected by both test types. These results contradict the accepted belief that unit and integration tests detect different types of defects in the software. This raises the question, if those definitions still fit to modern development contexts or if we should apply another distinction criterion to differentiate between unit and integration tests. It might not be a good idea to distinguish tests based on their properties, but on their usage as Ammann et al. [19] suggest. They reason that "most of the literature emphasizes these levels in terms of **when** they are applied, a more important distinction is on the **types of faults** that we are looking for." [19]. Hence, Ammann et al. [19] suggest that tests should be categorized according to their purpose or usage (i.e., types of defects targeted) and not according to their properties or when they are applied. While they use a different terminology, the core of the separation between unit and integration tests is similar to the definitions of the IEEE. This paper provides empirical data that it might not make sense to differentiate unit and integration tests in the way we nowadays do. Instead, we should create new definitions that fit to the modern software

development contexts and should follow the proposal of Ammann et al. [19] to separate tests based on which types of defects are targeted. Hence, we would need to revise current valid definitions that we teach in academia and industrial certifications. Our study results can provide valuable insights that could help with this revision.

**Practice:** Software development should always be accompanied by testing the developed parts. Our results show, that there seems to be no need to focus on the type of test developed, at least if we only consider the effectiveness of tests and classify them according to the IEEE definition. Nevertheless, other influences could play a role, e.g., if a test should serve as documentation [68]. Our results highlight, that there are differences from project to project. Hence, it seems that the context in which a project is developed has an influence on the effectiveness of the test types, as for some projects unit tests are more effective than integration tests, while for others it is vice versa. However, it seems that it is more important that developers "just create tests" instead of caring about mocking classes or which test type is developed.

To come back to the cited discussion from our introduction, where developers discuss if they should change their testing habit and focus more on integration than unit testing: our results show that this change would not have a negative impact, as both test types are equally effective. Hence, if developers argue that integration tests are more realistic and valuable in their daily developer life, our results do not argue against this practice. However, we also found that some defects were only detected by either unit or integration tests. Moreover, our results highlight that mutants that are hard to kill are (on average) more effectively killed by unit tests, as the results with the **DISJ** data set depict. Hence, our results show that it is still favorable to test on both test levels.

## 5. Threats to Validity

In this section, we discuss the external, internal, and construct validity of our study together with the validation procedures that we have taken to counter or measure those threats.

### 5.1. Construct Validity

Threats to this type of validity are concerned with the degree to which our analysis really measures what we intended it to measure. While we carefully tested our tools and scripts via manually written tests and manually curated samples of data, there can still be defects, which can influence our results.

For the collection of the mutation data we make use of PIT [52], which is a mature mutation testing framework that is often used in research, e.g. [53, 55]. Hence, it is less likely that the mutation testing is not working correctly. Every problem that occurred during the mutation testing was manually inspected. We found that PIT sometimes reported that a test was not finishing without a failure, but this only occurred for 35 out of 36435 test cases. Furthermore, PIT is designed for mutation testing on unit level. While PIT offers some mutation operators that are also offered by tools that are used for integration mutation testing (e.g. [57, 58]), this still could have an influence on our results.

The threat of mutant redundancy is substantial for our kind of analysis. We tried to reduce this thread by including the disjoint mutant set into our analysis process, as proposed by Kintis et al. [53]. Nevertheless, this is an approximation and it might not be the most fitting set of mutants.

Another threat is our choice of the defect classification scheme. A different classification scheme might produce different results for our research question. But, we have chosen the scheme by Zhao et al. [37] as it is close to the code, which is needed to classify integrated mutations, and it provides a good overview of different defect types.

Moreover, our approach to create a defect classification for an integrated mutation can be flawed for some mutations. For example, if a mutation is integrated in a line which has several statements on it so that the re-integration

of the defect in the correct source code could be done in the wrong statement. To measure this threat, we manually checked a sample of our defect type classifications. None of the defect type classifications had the problem mentioned above.

We reuse the build file of the projects that we analyze to execute all of their tests. Hence, it can be the case that some tests are filtered out, e.g., because they should only be executed via an CI system as they are very long running. It can be the case that these filtered tests are always tests from the same type (e.g., integration tests, as they often have a higher execution time than unit tests). Nevertheless, we assess this threat by counting the number of tests that are excluded from execution. Overall, 67 tests are excluded. However, this number is rather low in comparison to the overall number of executed tests (i.e. 38782).

*5.2. Internal Validity*

Internal validity threats are concerned with the ability to draw conclusions from the relation between causes and effects. While we tried to create an isolated environment where we have control over influencing variables (e.g., by using mutation analysis to integrate defects into a software) it can be the case that the defect detection capabilities are also influenced by other variables. We countered the influences that we know of, e.g., by normalizing the results.

The current version of our framework can not differentiate integration from system tests. Hence, it is possible that a system test is classified as an integration test. To mitigate this issue we included only projects that are libraries or frameworks.

The statistical tests used in this paper rely on the accurate implementations of the algorithms in the external library used. To ease this threat we are using a well known public library, namely SciPy [69], for these algorithms.

*5.3. External Validity*

Threats to this type of validity are concerned with the ability to generalize our results. We had a look at Java projects only. Hence, the results can be

different for projects written in other languages. Although we analyzed a larger sample of projects than most other related work, the results can vary if other projects are chosen. Especially, as our set of projects is limited to open-source projects, even if some of the projects are developed by companies in an open-source manner (e.g., *google-gson*). Furthermore, we only selected libraries or frameworks which could potentially influence our results. Hence, our results could be different for other project types (e.g., applications). But, our approach needs a compilable release of a project, which is often not given as Tufano et al. [70] highlight in their paper. This complicates a fully automatic analysis with more data. Nevertheless, the replication of this work using other projects (and other programming languages) is required in order to reach a more general conclusion. Hence, we added a replication kit that includes all data and programs of this paper to support such conceptual replications.

## 6. Conclusion and Future Work

In this paper, we reported an empirical study that was conducted on 17 open-source Java projects. The goal of this study was to investigate, if the existing standard definitions of unit and integration tests are still valid in modern software development contexts. We created two different data sets to pursue this goal: mutation testing data representing a large variety of defects that could be introduced in the program (627507 unique mutations) and the disjoint mutation testing data set that represents only "hard to kill" mutants. Additionally, we classified all created mutations into different defect types to evaluate if there is a defect type that is more effectively detected by a certain test type. Overall we collected and analyzed the defect detection capabilities of 38782 test cases.

Our results highlight, that there are no significant differences in the defect detection capabilities of unit and integration tests in either of our two data sets. However, we also found that there are some defects that can only be detected by one test type (i.e., either unit or integration tests). Furthermore, we found that we can not state that one defect type can be more effectively

found by a unit or an integration test, as our tests found no statistical significant differences. Hence, it seems that the current standard definitions do not fit to modern software development contexts anymore, as there should be a difference between those test types. Therefore, our results questions if the division between unit and integration tests is reasonable for modern systems, like we develop nowadays. These results suggest that we should create usage-based definitions instead of property-based definitions of unit and integration tests so that they fit to modern software development contexts.

Our future work includes but is not limited to the use of more projects with different programming languages for the performed and additional analyses. As additional analysis we plan a qualitative study on the defects and test cases that we used in this study. This could help us to understand the differences between unit and integration tests and their (not) detected defects. Furthermore, we would like to perform a developer study on unit and integration testing practices to get feedback from developers how they use unit and integration tests in their daily work. In connection to this study, it would be interesting to assess the usage of different test types in different development phases or development models, e.g., by comparing the usage of unit and integration tests in an classical and agile development environment.

[1] B. Van Rompaey, B. Du Bois, S. Demeyer, M. Rieger, On the detection of test smells: A metrics-based approach for general fixture and eager test, IEEE Transactions on Software Engineering 33 (12) (2007) 800–817.

[2] G. Fraser, A. Zeller, Generating parameterized unit tests, in: Proceedings of the 20th ACM SIGSOFT International Symposium on Software Testing and Analysis, ACM, 2011, pp. 364–374.

[3] A. Gambi, S. Kappler, J. Lampel, A. Zeller, Cut: automatic unit testing in the cloud, in: Proceedings of the 26th ACM SIGSOFT International Symposium on Software Testing and Analysis, ACM, 2017, pp. 364–367.

[4] R. Lachmann, S. Lity, M. Al-Hajjaji, F. Fürchtegott, I. Schaefer, Fine-grained test case prioritization for integration testing of delta-oriented software product lines, in: Proceedings of the 7th International Workshop on Feature-Oriented Software Development, ACM, 2016, pp. 1–10.

[5] D. Xu, W. Xu, M. Tu, N. Shen, W. Chu, C.-H. Chang, Automated integration testing using logical contracts, IEEE Transactions on Reliability 65 (3) (2016) 1205–1222.

[6] D. Holling, A. Hofbauer, A. Pretschner, M. Gemmar, Profiting from unit tests for integration testing, in: IEEE International Conference on Software Testing, Verification and Validation (ICST), IEEE, 2016, pp. 353–363.

[7] B. W. Boehm, Verifying and validating software requirements and design specifications, IEEE Software 1 (1) (1984) 75.

[8] W. Royce, Software project management, Pearson Education India, 1999.

[9] IEEE, Systems and software engineering – vocabulary, ISO/IEC/IEEE 24765:2010(E) (2010) 1–418doi:10.1109/IEEESTD.2010.5733835.

[10] International Software Testing Qualification Board, International Software Testing Qualification Board Glossary, http://www.astqb.org/glossary/search/unit, [accessed 30-November-2018].

[11] JUnit Team, JUnit Homepage, http://junit.org/junit5/, [accessed 30-November-2018] (2017).

[12] Travis CI GmbH, Travis CI Homepage, https://travis-ci.org/, [accessed 30-November-2018] (2017).

[13] Jenkins Contributors, Jenkins, https://jenkins.io/, [accessed 30-November-2018] (2018).

[14] K. C. Dodds, Write tests. Not too many. Mostly integration., `https://blog.kentcdodds.com/write-tests-not-too-many-mostly-integration-5e8c7fff591c`, [accessed 30-November-2018] (2017).

[15] J. O. Coplien, Seque, `http://rbcs-us.com/documents/Segue.pdf`, [accessed 30-November-2018] (2014).

[16] J. O. Coplien, Why Most Unit Testing is Waste, `http://rbcs-us.com/documents/Why-Most-Unit-Testing-is-Waste.pdf`, [accessed 30-November-2018] (2014).

[17] F. Trautsch, J. Grabowski, Are there any unit tests? an empirical study on unit testing in open source python projects, in: IEEE International Conference on Software Testing, Verification and Validation (ICST), IEEE, 2017, pp. 207–218.

[18] A. Spillner, T. Linz, H. Schaefer, Software testing foundations: a study guide for the certified tester exam, Rocky Nook, Inc., 2014.

[19] P. Ammann, J. Offutt, Introduction to software testing, Cambridge University Press, 2016.

[20] G. J. Myers, C. Sandler, T. Badgett, The art of software testing, John Wiley & Sons, 2011.

[21] G. Orellana, G. Laghari, A. Murgia, S. Demeyer, On the differences between unit and integration testing in the travistorrent dataset, in: Proceedings of the 14th International Conference on Mining Software Repositories, IEEE Press, 2017, pp. 451–454.

[22] Apache Software Foundation, Maven Surefire Plugin Website, `http://maven.apache.org/surefire/maven-surefire-plugin/`, [accessed 30-November-2018] (2017).

[23] Apache Software Foundation, Maven FailSafe Plugin Website, `http://maven.apache.org/surefire/maven-failsafe-plugin/`, [accessed 30-November-2018] (2017).

[24] T. Kanstrén, Towards a deeper understanding of test coverage, Journal of Software: Evolution and Process 20 (1) (2008) 59–76.

[25] M. Beller, G. Gousios, A. Zaidman, Travistorrent: Synthesizing travis ci and github for full-stack research on continuous integration, in: Proceedings of the 14th International Conference on Mining Software Repositories, 2017.

[26] M. Papadakis, M. Kintis, J. Zhang, Y. Jia, Y. Le Traon, M. Harman, Mutation testing advances: an analysis and survey, Advances in Computers.

[27] M. Papadakis, Y. Jia, M. Harman, Y. Le Traon, Trivial compiler equivalence: A large scale empirical study of a simple, fast and effective equivalent mutant detection technique, in: Proceedings of the 37th International Conference on Software Engineering-Volume 1, IEEE Press, 2015, pp. 936–946.

[28] D. Schuler, A. Zeller, Covering and uncovering equivalent mutants, Software Testing, Verification and Reliability 23 (5) (2013) 353–374.

[29] M. Daran, P. Thévenod-Fosse, Software error analysis: A real case study involving real faults and mutations, in: ACM SIGSOFT Software Engineering Notes, Vol. 21, ACM, 1996, pp. 158–171.

[30] J. H. Andrews, L. C. Briand, Y. Labiche, Is mutation an appropriate tool for testing experiments?, in: Proceedings of the 27th International Conference on Software Engineering, ACM, 2005, pp. 402–411.

[31] J. H. Andrews, L. C. Briand, Y. Labiche, A. S. Namin, Using mutation analysis for assessing and comparing testing coverage criteria, IEEE Transactions on Software Engineering 32 (8) (2006) 608–624.

[32] R. Just, D. Jalali, L. Inozemtseva, M. D. Ernst, R. Holmes, G. Fraser, Are mutants a valid substitute for real faults in software testing?, in: Proceed-

ings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering, ACM, 2014, pp. 654–665.

[33] A. S. Namin, S. Kakarla, The use of mutation in testing experiments and its sensitivity to external threats, in: Proceedings of the 2011 International Symposium on Software Testing and Analysis, ACM, 2011, pp. 342–352.

[34] T. T. Chekam, M. Papadakis, Y. L. Traon, M. Harman, An empirical study on mutation, statement and branch coverage fault revelation that avoids the unreliable clean program assumption, in: Proceedings of the 39th International Conference on Software Engineering, IEEE Press, 2017, pp. 597–608.

[35] M. Papadakis, D. Shin, S. Yoo, D.-H. Bae, Are mutation scores correlated with real fault detection? a large scale empirical study on the relationship between mutants and real faults, in: 40th International Conference on Software Engineering, May 27-3 June 2018, Gothenburg, Sweden, 2018.

[36] J. H. Hayes, Building a requirement fault taxonomy: Experiences from a nasa verification and validation research project, in: 14th International Symposium on Software Reliability Engineering, IEEE, 2003, pp. 49–59.

[37] Y. Zhao, H. Leung, Y. Yang, Y. Zhou, B. Xu, Towards an understanding of change types in bug fixing code, Information and Software Technology 86 (2017) 37–53.

[38] X. Xia, D. Lo, X. Wang, B. Zhou, Automatic defect categorization based on fault triggering conditions, in: Proceedings of the 19th International Conference on Engineering of Complex Computer Systems (ICECCS), IEEE, 2014, pp. 39–48.

[39] L. Tan, C. Liu, Z. Li, X. Wang, Y. Zhou, C. Zhai, Bug characteristics in open source software, Empirical Software Engineering 19 (6) (2014) 1665–1705.

[40] R. Chillarege, I. S. Bhandari, J. K. Chaar, M. J. Halliday, D. S. Moebus, B. K. Ray, M.-Y. Wong, Orthogonal defect classification-a concept for in-process measurements, IEEE Transactions on Software Engineering 18 (11) (1992) 943–956.

[41] A. J. Offutt, J. H. Hayes, A semantic model of program faults, in: ACM SIGSOFT Software Engineering Notes, Vol. 21, ACM, 1996, pp. 195–200.

[42] J. H. Hayes, C. I. Raphael, V. K. Surisetty, A. Andrews, Fault links: exploring the relationship between module and fault types, in: European Dependable Computing Conference, Springer, 2005, pp. 415–434.

[43] Q. Tu, et al., Evolution in open source software: A case study, in: Proceedings of the Internation Conference on Software Maintenance, IEEE, 2000, pp. 131–142.

[44] Apache Software Foundation, Maven Project Homepage, `https://maven.apache.org/`, [accessed 30-November-2018] (2017).

[45] C. Beust, TestNG Documentation, `http://testng.org/doc/`, [accessed 30-November-2018] (2015).

[46] H. Borges, A. Hora, M. T. Valente, [dataset] improved list of popular github repositories, `https://doi.org/10.5281/zenodo.804473`, [accessed 30-November-2018] (2017).

[47] MVN Repository, MVN Repository - Most Popular, `https://mvnrepository.com/popular`, [accessed 30-November-2018] (2018).

[48] F. Trautsch, S. Herbold, P. Makedonski, J. Grabowski, Adressing problems with external validity of repository mining studies through a smart data platform, in: Proceedings of the 13th International Conference on Mining Software Repositories, ACM, 2016, pp. 97–108.

[49] N. C. Borle, M. Feghhi, E. Stroulia, R. Greiner, A. Hindle, Analyzing the effects of test driven development in github, Empirical Software Engineering (2017) 1–28.

[50] Oracle, What Is a Package?, `https://docs.oracle.com/javase/tutorial/java/concepts/package.html`, [accessed 30-November-2018] (2017).

[51] Apache Software Foundation, commons-io GitHub, `https://github.com/apache/commons-io`, [accessed 30-November-2018] (2018).

[52] H. Coles, T. Laurent, C. Henard, M. Papadakis, A. Ventresque, Pit: a practical mutation testing tool for java, in: Proceedings of the 25th International Symposium on Software Testing and Analysis, ACM, 2016, pp. 449–452.

[53] M. Kintis, M. Papadakis, A. Papadopoulos, E. Valvis, N. Malevris, Y. Le Traon, How effective are mutation testing tools? an empirical analysis of java mutation testing tools with manual analysis and real faults, Empirical Software Engineering (2017) 1–38.

[54] H. Coles, PIT Project Homepage, `http://pitest.org/`, [accessed 30-November-2018] (2017).

[55] T. Laurent, M. Papadakis, M. Kintis, C. Henard, Y. Le Traon, A. Ventresque, Assessing and improving the mutation testing practice of pit, in: IEEE International Conference on Software Testing, Verification and Validation (ICST), IEEE, 2017, pp. 430–435.

[56] M. Papadakis, C. Henard, M. Harman, Y. Jia, Y. Le Traon, Threats to the validity of mutation-based test assessment, in: Proceedings of the 25th International Symposium on Software Testing and Analysis, ACM, 2016, pp. 354–365.

[57] M. E. Delamaro, J. Maidonado, A. P. Mathur, Interface mutation: An approach for integration testing, IEEE Transactions on Software Engineering 27 (3) (2001) 228–247.

[58] M. Grechanik, G. Devanla, Mutation integration testing, in: IEEE International Conference on Software Quality, Reliability and Security (QRS), IEEE, 2016, pp. 353–364.

[59] F. Trautsch, BugFixClassifier GitHub, `https://github.com/ftrautsch/BugFixClassifier`, [accessed 30-November-2018] (2018).

[60] B. Fluri, H. C. Gall, Classifying change types for qualifying change couplings, in: 14th IEEE International Conference on Program Comprehension, IEEE, 2006, pp. 35–45.

[61] R. Wang, S. W. Lagakos, J. H. Ware, D. J. Hunter, J. M. Drazen, Statistics in medicinereporting of subgroup analyses in clinical trials, New England Journal of Medicine 357 (21) (2007) 2189–2194.

[62] S. S. Shapiro, M. B. Wilk, An analysis of variance test for normality (complete samples), Biometrika 52 (3/4) (1965) 591–611.

[63] I. Olkin, Contributions to probability and statistics: essays in honor of Harold Hotelling, Stanford University Press, 1960.

[64] Student, The probable error of a mean, Biometrika 6 (1) (1908) 1–25. `doi:10.1093/biomet/6.1.1`.

[65] H. B. Mann, D. R. Whitney, On a test of whether one of two random variables is stochastically larger than the other, The annals of mathematical statistics (1947) 50–60.

[66] M. Aickin, H. Gensler, Adjusting for multiple testing when reporting research results: the bonferroni vs holm methods., American journal of public health 86 (5) (1996) 726–728.

[67] F. Trautsch, S. Herbold, J. Grabowski, [dataset] Replication Kit, `https://doi.org/10.5281/zenodo.2267946`, [accessed 30-November-2018] (2018).

[68] S. Demeyer, S. Ducasse, O. Nierstrasz, Object-oriented reengineering patterns, Elsevier, 2002.

[69] SciPy developers, SciPy Website, `https://www.scipy.org/`, [accessed 30-November-2018] (2017).

[70] M. Tufano, F. Palomba, G. Bavota, M. Di Penta, R. Oliveto, A. De Lucia, D. Poshyvanyk, There and back again: Can you compile that snapshot?, Journal of Software: Evolution and Process 29 (4).

## Appendix A. Change Type Mapping

| Change Type | Defect Class |
| --- | --- |
| ADDING_ATTRIBUTE_MODIFIABILITY | Data |
| ADDING_CLASS_DERIVABILITY | Interface |
| ADDING_METHOD_OVERRIDABILITY | Interface |
| ADDITIONAL_CLASS | Interface |
| ADDITIONAL_OBJECT_STATE | Data |
| ALTERNATIVE_PART_DELETE | Logic/Control |
| ALTERNATIVE_PART_INSERT | Logic/Control |
| ATTRIBUTE_RENAMING | Data |
| ATTRIBUTE_TYPE_CHANGE | Data |
| CLASS_RENAMING | Interface |
| COMMENT_DELETE | Other |
| COMMENT_INSERT | Other |
| COMMENT_MOVE | Other |
| COMMENT_UPDATE | Other |
| CONDITION_EXPRESSION_CHANGE | Logic/Control |
| DECREASING_ACCESSIBILITY_CHANGE | Interface |
| DOC_DELETE | Other |
| DOC_INSERT | Other |

| | |
|---|---|
| DOC_UPDATE | Other |
| INCREASING_ACCESSIBILITY_CHANGE | Interface |
| METHOD_RENAMING | Interface |
| PARAMETER_DELETE | Interface |
| PARAMETER_INSERT | Interface |
| PARAMETER_ORDERING_CHANGE | Interface |
| PARAMETER_RENAMING | Interface |
| PARAMETER_TYPE_CHANGE | Interface |
| PARENT_CLASS_CHANGE | Interface |
| PARENT_CLASS_DELETE | Interface |
| PARENT_CLASS_INSERT | Interface |
| PARENT_INTERFACE_CHANGE | Interface |
| PARENT_INTERFACE_DELETE | Interface |
| PARENT_INTERFACE_INSERT | Interface |
| REMOVED_CLASS | Interface |
| REMOVED_OBJECT_STATE | Data |
| REMOVING_ATTRIBUTE_MODIFIABILITY | Data |
| REMOVING_CLASS_DERIVABILITY | Interface |
| REMOVING_METHOD_OVERRIDABILITY | Interface |
| RETURN_TYPE_CHANGE | Interface |
| RETURN_TYPE_DELETE | Interface |
| RETURN_TYPE_INSERT | Interface |

**Table A.7** Mapping of the change types by Fluri et al. [60] that can be *directly* mapped onto the defect classes by Zhao et al. [37].

| Condition | Defect Class |
|---|---|
| $CT \in \{\texttt{STATEMENT\_*}\} \wedge$<br>$CE \in \{\texttt{ASSIGNMENT, POSTFIX\_EXPRESSION, PREFIX\_EXPRESSION}\} \wedge$<br>$PE \notin \{\texttt{FOR\_INCR}\}$ | Computation |
| $CT \in \{\texttt{STATEMENT\_*}\} \wedge$<br>$CE \in \{\texttt{VARIABLE\_DECLARATION\_STATEMENT}\} \wedge$<br>$PE \notin \{\texttt{FOR\_INIT}\}$ | Data |
| $CT \in \{\texttt{UNCLASSIFIED\_CHANGE}\} \wedge$<br>$CE \in \{\texttt{MODIFIER}\}$ | Data |
| $CT \in \{\texttt{STATEMENT\_*}\} \wedge$<br>$CE \in \{\texttt{METHOD\_INVOCATION, CONSTRUCTOR\_INVOCATION,}$<br>$\texttt{SYNCHRONIZED\_STATEMENT, CLASS\_INSTANCE\_CREATION}\}$ | Interface |
| $CT \in \{\texttt{ADDING\_FUNCTIONALITY, REMOVING\_FUNCTIONALITY}\} \wedge$<br>$CE \in \{\texttt{METHOD}\}$ | Interface |
| $CT \in \{\texttt{UNCLASSIFIED\_CHANGE}\} \wedge$<br>$CE \in \{\texttt{TYPE\_PARAMETER}\}$ | Interface |
| $CT \in \{\texttt{STATEMENT\_*}\} \wedge$<br>$CE \in \{\texttt{IF\_STATEMENT, FOREACH\_STATEMENT, CONTINUE\_STATEMENT,}$<br>$\texttt{RETURN\_STATEMENT, THROW\_STATEMENT, SWITCH\_CASE,}$<br>$\texttt{SWITCH\_STATEMENT, BREAK\_STATEMENT, CATCH\_CLAUSE,}$<br>$\texttt{TRY\_STATEMENT, FOR\_STATEMENT, WHILE\_STATEMENT, DO\_STATEMENT,}$<br>$\texttt{LABELED\_STATEMENT}\}$ | Logic/Control |
| $CT \in \{\texttt{STATEMENT\_*}\} \wedge$<br>$CE \in \{\texttt{ASSIGNMENT, POSTFIX\_EXPRESSION, PREFIX\_EXPRESSION}\} \wedge$<br>$PE \in \{\texttt{FOR\_INCR}\}$ | Logic/Control |
| $CT \in \{\texttt{STATEMENT\_*}\} \wedge$<br>$CE \in \{\texttt{VARIABLE\_DECLARATION\_STATEMENT}\} \wedge$<br>$PE \in \{\texttt{FOR\_INIT}\}$ | Logic/Control |
| $CT \in \{\texttt{STATEMENT\_*}\} \wedge$<br>$CE \in \{\texttt{ASSERT\_STATEMENT}\}$ | Other |

**Table A.8** Mapping of the change types (CT) by Fluri et al. [60], where the changed entity (CE) and/or the parent entity (PE) needs to be taken into account to map a change onto the defect classes by Zhao et al. [37]. The term STATEMENT_* includes the general change types, i.e., `STATEMENT_UPDATE`, `STATEMENT_INSERT`, `STATEMENT_DELETE`, `STATEMENT_PARENT_CHANGE`, `STATEMENT_ORDERING_CHANGE`.