

Measuring Software Testability Modulo Test Quality

Valerio Terragni

USI Università della Svizzera italiana
Switzerland
valerio.terragni@usi.ch

Pasquale Salza

University of Zurich
Switzerland
salza@ifi.uzh.ch

Mauro Pezzè

USI Università della Svizzera italiana
Schaffhausen Institute of Technology
Switzerland
mauro.pezze@usi.ch

ABSTRACT

Comprehending the degree to which software components support testing is important to accurately schedule testing activities, train developers, and plan effective refactoring actions. Software testability estimates such property by relating code characteristics to the test effort. The main studies of testability reported in the literature investigate the relation between class metrics and test effort in terms of the size and complexity of the associated test suites. They report a moderate correlation of some class metrics to test-effort metrics, but suffer from two main limitations: (i) the results hardly generalize due to the small empirical evidence (datasets with no more than eight software projects); and (ii) mostly ignore the quality of the tests. However, considering the quality of the tests is important. Indeed, a class may have a low test effort because the associated tests are of poor quality, and not because the class is easier to test. In this paper, we propose an approach to measure testability that normalizes the test effort with respect to the test quality, which we quantify in terms of code coverage and mutation score. We present the results of a set of experiments on a dataset of 9,861 JAVA classes, belonging to 1,186 open source projects, with around 1.5 million of lines of code overall. The results confirm that normalizing the test effort with respect to the test quality largely improves the correlation between class metrics and the test effort. Better correlations result in better prediction power and thus better prediction of the test effort.

CCS CONCEPTS

• **Software and its engineering** → **Designing software**; **Software testing and debugging**; *Software libraries and repositories*; *Software design tradeoffs*.

KEYWORDS

Software Testability, Test Effort, Test Quality, Software Metrics

ACM Reference Format:

Valerio Terragni, Pasquale Salza, and Mauro Pezzè. 2020. Measuring Software Testability Modulo Test Quality. In *28th International Conference on Program Comprehension (ICPC '20)*, October 5–6, 2020, Seoul, Republic of Korea. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3387904.3389273>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ICPC '20, October 5–6, 2020, Seoul, Republic of Korea

© 2020 Association for Computing Machinery.

ACM ISBN 978-1-4503-7958-8/20/05...\$15.00

<https://doi.org/10.1145/3387904.3389273>

1 INTRODUCTION

Software testing is an essential, labor-intensive and time-consuming activity of the software life cycle. Making testing easier is important for many software companies, as it lowers development costs while increasing the number of detected faults.

It is well understood that the effort of testing software systems depends on the artifacts under test, meaning that some software systems are easier to test than others [13, 22, 54]. Comprehending the relation between software artifacts and test effort is extremely important to control the cost of testing and improve the accuracy of test plans. *Software Testability* [22] captures the impact of software artifacts on testing by estimating *the degree to which a software system or component under test supports its own testing*¹. A software system with a high degree of testability results in a low test effort.

In their recent comprehensive literature review of 208 papers on software testability, Garousi et al. [22] observe that measuring and predicting the testability is the topic that received the most attention [2–4, 12, 13, 51]. The general idea is to measure (predict) the test effort of software systems from structural metrics of the software that are available before designing the test cases [51]. Early predicting the test effort can help developers to (i) early identify software components that require more test effort, on which developers have to focus to ensure software quality, (ii) plan testing activities and optimally allocate resources, and (iii) recognize refactoring opportunities to reduce the test effort.

Most studies on measuring and predicting testability investigate on the relation between class-level metrics in object oriented systems, for instance Chidamber and Kemerer (C&K) [14], and the cost of writing test cases (the test effort) [2–4, 12, 13, 51]. These studies approximate the test effort with the size and complexity of the test suites, for instance the number of tests and assertions in the test class associated to the class under test. They provide some evidence of the existence of a correlation between the class-level metrics and the test effort, but suffer from two limitations: (i) the data sets are of small size, and (ii) mostly ignore the quality of the test suites.

Small sample size. Previous studies involved at most eight software projects [22]. Such a small number of analyzed projects does not guarantee the generalizability of the results: specific development styles, frameworks, and practices can influence the correlation results and produce different results for different projects [3].

Ignoring the test quality. Previous studies measured the test effort in terms of the size of the test classes, while mostly ignoring

¹The literature proposes many definitions of software testability [22]. In this paper, we refer to both the IEEE 610.12-1990 and ISO/IEC 9126 standards [54] that define testability in similar ways: IEEE: “the degree to which a system or a component facilitates the establishment of test criteria and the performance of tests to determine whether those criteria have been met.” ISO: “attributes of software that bear on the effort needed to validate the software product.”

the quality of the tests. Lacking a quality assessment of the tests leads to imprecise correlation results: classes with comparable test effort but different test quality should not have the same degree of testability. A class may have a low test effort because the associated tests are of poor quality and not because the class is easier to test. Bruntink and van Deursen’s study is the only work that acknowledges the test quality when comparing with the test effort [13]. They ensured that the analyzed software systems had test suites of similar quality, measured in terms of line coverage. However, their correlation study involves only five software systems [13].

In this paper, we propose a new approach to measure the testability of object-oriented classes. *Our approach normalizes the test effort of a class with respect to the quality of its tests*, which we quantify with code coverage and mutation score. This enables the correlation analyses and prediction models of an arbitrary large number of heterogeneous software systems, implemented with different test-quality criteria.

We empirically investigated our approach with 28 metrics to characterize the class properties, six metrics to measure the test effort, and three metrics to quantify the test quality. We analyzed 9,861 pairs of JAVA classes and corresponding JUNIT test classes collected from 1,186 open source projects on GITHUB. We computed the *Spearman’s correlation coefficient* (ρ) [27] for all 168 pair-wise combinations of class and test-effort metrics, before and after the normalization with the test-quality metrics.

The results confirm that some class metrics correlates with test-effort metrics, and indicate that normalizing test-effort metrics with test-quality metrics drastically improves the correlation (up to 74%). Better correlation leads to better prediction power, and thus better prediction of test effort [51]. On the one hand, we use the test-quality metrics to normalize the test effort for the correlation analysis, allowing to fairly compare classes belonging to projects of different test qualities. On the other hand, the test-quality metrics can also be used as a target variable for prediction purposes. Indeed, if the purpose is to predict the test effort before writing the tests, a target value for test quality can be used in a preprocessing step to normalize the dataset used for training a prediction model. For instance, one might want to predict the test effort required to write tests having a target mutation score of 80%. The data used to train the model can be normalized according to that value, to build a model that predicts the test effort for the targeted mutation score.

Our empirical study concludes that (i) normalizing test-effort metrics by mutation score achieves the best correlation improvement, and (ii) the object-oriented design properties that most influence testability are: size, complexity, coupling, and cohesion.

This paper contributes to a better comprehension of software testability by:

- presenting the by-far largest study on the correlation of class and test-effort metrics in terms of analyzed metrics, classes and projects;
- extending the testability measurements by normalizing the test effort, with respect to the quality of the test suites;
- showing that the proposed normalization improves the correlation between class metrics and test effort.
- giving important insights on software testability that confirm some of the findings of previous studies as well as uncover

previously unknown correlations between object-oriented design properties and test effort;

- publicly releasing our dataset for further studies².

The paper is organized as follows. Section 2 presents the objective of this study, introduces the considered metrics, and motivates and explains our normalization procedure. Section 3 presents the results, which address the main research questions that validate the hypothesis that “normalizing by the test quality” achieves better correlation than not using the normalization. Section 4 discusses the related work. Section 5 summarizes the main results presented in the paper.

2 OBJECTIVE AND METHODOLOGY

This paper investigates the relationships between the object-oriented metrics of classes and the test effort of designing unit test cases for such classes. We produced statistically significant and general results for JAVA software systems, by conducting an experimental study on a large set of heterogeneous JAVA software systems of different size and category. Because heterogeneous projects are likely to have different test-quality criteria, we introduce a novel normalization procedure that homogenizes the values of test-effort metrics according to the values of test-quality metrics.

This section contextualizes the scope of the study (Section 2.1), presents the considered metrics of class, test effort and test quality (Section 2.2, Section 2.3, and Section 2.4, respectively), and introduces the new normalization procedure (Section 2.5).

2.1 Testability of Object-Oriented Programs

We target systems designed with the **Object Oriented (OO)** programming paradigm [45], which is based on the concept of “objects” (instances of classes) that can contain both data (object fields) and code (methods). In particular, we consider systems written in the JAVA language.

Testing OO programs is often performed at three different levels [55]: unit, integration, and system. *Unit testing* tests in isolation small portions of programs called units, for instance methods or classes. The goal of unit testing is to isolate each part of the program and show that individual parts are correct. *Integration testing* tests the interaction of multiple units. *System testing* tests a complete and integrated software system.

When dealing with software testability, unit testing is the most useful testing level because one can apply testability analysis early in the development life-cycle [13]. In line with the work presented in the literature [22], we study software testability at **unit level**, and more specifically at **class level**.

Working at class level testing has two practical advantages. First, we can leverage several OO metrics defined at class level [5, 14, 47]. Second, we can take advantage of popular naming conventions to identify the test class associated with a given class [13]. In fact, a common software development practice in OO programming languages, such as JAVA, is to create a dedicated class for each tested class following well-defined naming conventions [21, 39].

²The dataset and source code for the analysis we performed in this paper are shared at the address <https://doi.org/10.5281/zenodo.3740499>.

Table 1: Class metrics

Design Property	Name	Description	Reference
Size	Lines of Code (LOC)	Number of non-blank lines including comments and JavaDoc	-
	Number of Bytecode Instructions (NBI)	Number of bytecode instructions in the compiled .class file	-
	Lines of Comment (LOCCOM)	Number of lines of comment, excluding any end-of-line comments	-
	Number of Public Methods (NPM)	Number of methods in a class that are declared public	Goyal and Joshi [23]
	Number of STatic Method (NSTAM)	Number of methods in a class that are declared static	-
	Number of Fields (NOF)	Number of fields (attributes) in the class	-
	Number of STatic Fields (NSTAF)	Number of static fields (attributes) in the class	-
	Number of Method Calls (NMC)	Number of method invocations	-
	Number of Method Calls Internal (NMCI)	Number of method invocations of methods defined in the class	-
	Number of Method Calls External (NMCE)	Number of method invocations of methods defined in other classes	-
Complexity	Weighted Methods per Class (WMC)	Sum of the Cyclomatic Complexity [38] of all methods in the class	Chidamber and Kemerer [14]
	Average Method Complexity (AMC)	Average of the Cyclomatic Complexity [38] of all methods in the class	Tang, Kao and Chen [47]
	Response For a Class (RFC)	Number of methods that response to a message from the class itself	Chidamber and Kemerer [14]
Inheritance	Depth of Inheritance Tree (DIT)	Number of super classes	Chidamber and Kemerer [14]
	Number of Children (NOC)	Number of immediate sub-classes subordinated to a class in the class hierarchy	Chidamber and Kemerer [14]
	Measure of Functional Abstraction (MFA)	Ratio of the number of methods inherited by the class to the number of methods	Goyal and Joshi [23]
Coupling	Coupling Between Object classes (CBO)	Number of other classes that a class is coupled to	Chidamber and Kemerer [14]
	Inheritance Coupling (IC)	Number of parent classes to which a given class is coupled	Tang, Kao and Chen [47]
	Coupling Between Methods (CBM)	Number of redefined methods to which all the inherited methods are coupled	Tang, Kao and Chen [47]
	Afferent Coupling (Ca)	Measure of how many other classes use the specific class	Martin [36]
	Efferent Coupling (Ce)	Measure of how many other classes is used by the specific class	Martin [36]
Cohesion	Lack of Cohesion in Methods (LCOM)	Diff. between the number of method pairs without and with common variables	Chidamber and Kemerer [14]
	Lack of Cohesion Of Methods (LCOM3)	Revised version of LCOM	Henderson-Sellers [28]
	Cohesion Among Methods in class (CAM)	Represents the relatedness among methods of a class	Goyal and Joshi [23]
Encapsulation	Data Access Metrics (DAM)	Ratio of the number of private fields to the total number of fields	Goyal and Joshi [23]
	Number of PRivate Fields (NPRIF)	Number of private fields (attributes) of the class	-
	Number of PRivate Methods (NPRIM)	Number of private methods of the class	-
	Number of PROtected Methods (NPROM)	Number of protected methods of the class	-

2.2 Class Metrics

Table 1 shows the 28 class metrics that we considered in this study. We tried to be as inclusive as possible when selecting the metrics, and avoided to limit the selection to metrics known to be correlated with testability [22]. This is because we also aimed to find unknown correlations that may arise with a large study. We considered well-known object oriented metrics: Chidamber and Kemerer (C&K) [14] and the Tang, Kao and Chen (TKC) [47] metrics. We enriched this already large set of metrics with additional metrics that may correlate with testability. Table 1 groups the 28 metrics based on the design property that each metric characterizes: *size*, *complexity*, *inheritance*, *coupling*, *cohesion*, and *encapsulation*.

Size. Size metrics includes standard ones such as Lines Of Code (LOC), metrics about the number of (static) methods and fields in the class (Number of Public Methods (NPM), Number of STatic Method (NSTAM), Number of Fields (NOF), and Number of STatic Fields (NSTAF)) and metrics about the number of internal and external method calls (Number of Method Calls (NMC), Number of Method Calls Internal (NMCI), Number of Method Calls External (NMCE)). In this paper we propose Number of Bytecode Instructions (NBI) as a new metric defined as the number of bytecode instructions in the compiled .class file. NBI can be more informative than the classic LOC metric, because single lines of code in JAVA can correspond to simple statements (for instance, the variable assignment) or complex statements (such as the lambda expressions), which correspond to few or many numbers of bytecode instructions, respectively. Thus, the NBI metrics distinguishes classes with a similar number of LOCs but with different types of statements (complex and simple).

Complexity. Complexity metrics include Weighted Methods per Class (WMC) and Average Method Complexity (AMC). Both metrics depend on the *cyclomatic complexity* metric [38], which is its number of the linearly-independent paths of a method [38]. Because the cyclomatic complexity is defined at the method level, WMC [14] and AMC [47] convert it to class-level by summing and averaging the cyclomatic complexities of all methods in the class, respectively.

Inheritance. Inheritance metrics capture the different aspects of the inheritance of a class. Depth of Inheritance Tree (DIT) is the number of super-classes, Number of Children (NOC) is the number of immediate sub-classes in the class hierarchy, and Measure of Functional Abstraction (MFA) is the ratio of the number of methods inherited by the class to the total number of methods in the class.

Coupling. Coupling metrics characterize the degree of interdependence between classes and methods. Classes that have a high (outgoing) efferent coupling use other parts of the system, increasing the possible execution paths [43]. Low coupling is often a sign of a well-structured software system and a good design.

Cohesion. Cohesion describes the binding of the elements within one method and within one object class, respectively. Low cohesion means that the class does a great variety of actions.

Encapsulation. Encapsulation metrics capture the degree of encapsulation of the classes. For instance, the number of private methods and fields (NPRIM and NPRIF), the ratio of the number of private fields to the total number of fields (DAM).

Table 2: Test-effort metrics

Name	Description	Reference
TEST – Lines Of Code (T-LOC)	Number of non-blank lines including comments and JavaDoc of the test class	Bruntink and van Deursen [13]
TEST – Number Of Tests (T-NOT)	Number of test cases in the test class (methods with the @Test annotation)	Bruntink and van Deursen [13]
TEST – Number Of Assertions (T-NOA)	Number of test assertions in the test class (invocations to org.junit.Assert)	Bruntink and van Deursen [13]
TEST – Number of Method Calls (T-NMC)	Number of method invocations in the test class	Toure et al. [49, 50]
TEST – Weighted Methods per Class (T-WMC)	Sum of the Cyclomatic Complexity [38] of all methods in the test class	Chidamber and Kemerer [14]
TEST – Average Method Complexity (T-AMC)	Average of the Cyclomatic Complexity [38] of all methods in the test class	Tang, Kao and Chen [47]

Table 3: Test-quality metrics of T_C with respect to its associated class C

Name	Description	Reference
Line coverage (L)	Ratio of source code lines in C that are executed by at least one test in T_C	–
Branch coverage (B)	Ratio of branches in C that are executed by at least one test in T_C	–
Mutation score (M)	Mutation score of a test suite T_C with respect to a class C is the ratio of mutants that are killed by T_C	De Millo et al. [20]

2.3 Test-Effort Metrics

Test-effort metrics measure the effort of testing classes in terms of the size and complexity of the associated test cases. We refer to test classes written in JUNIT, the most popular testing framework for JAVA [21, 39]. Let T_C denote the JUNIT **test class** of a class C . Each test method in T_C includes zero or more assertion oracles (boolean conditions) that predicate on the behavior of the class under test, and whose runtime values determine the pass or fail status of the test case. A test class may declare additional methods and inner classes to support the test executions, and such test code is called *test scaffolding*. Examples of scaffolding methods are those annotated with @Before and @After.

The effort of testing a class C is best quantified using the man-hours required to design and implement T_C [37]. However, collecting such information is difficult even for a few small software projects [37], and becomes unrealistic for many large software projects. As such, researchers often approximate the test effort with the size and complexity of the test class (*test-effort metrics*) [2, 3, 12, 13, 51], under the assumption that the size and complexity of a test class reflect the effort for designing it.

Table 2 shows the six test-effort metrics that we use to characterize test effort. The first four metrics TEST - Lines Of Code (T-LOC), TEST – Number Of Tests (T-NOT), TEST – Number Of Assertions (T-NOA), and TEST – Number of Method Calls (T-NMC) measure the size of the test class under different perspectives. T-LOC considers the size of the entire test class, and thus it includes scaffolding methods and inner classes. T-LOC also considers comment lines, since adding comments to the test code contributes to the cost of designing test cases [2]. T-NOT indicates the number of test cases regardless of their size. T-NOA is the amount of assertions contained in the test class, which might not be the same as T-NOT, because a test may have multiple assertions. T-NMC (called TINVOKE in Toure et al.’s paper [51]) measures the number of method calls in the test class, which is a proxy for the degree of dependency of the test cases [51]. According to Toure et al., T-NMC is particularly important to characterize the test effort of classes [50]. Intuitively, a test class with few dependencies is easier to execute than a test class with many dependences, because a test class with few dependencies can invoke the methods under test directly [51]. TEST –

Weighted Methods per Class (T-WMC) and TEST – Average Method Complexity (T-AMC) measure the cyclomatic complexity [38] of the test class, which well quantifies the test effort. A test class may contain test or scaffolding methods with many linearly independent execution paths. We assume that the higher the complexity of the test class the higher the effort required for writing and designing it.

2.4 Test-Quality Metrics

In this paper, we propose the usage of test-quality metrics to normalize test-effort metrics. Table 3 shows the three test-quality metrics that we considered: *line coverage*, *branch coverage*, and *mutation score*. These metrics are commonly used to approximate the quality of a test suite defined as the ability to reveal faults [54].

Coverage analysis [40] executes the test class T_C on an instrumented version of the class under test C to collect coverage information, and computes the percentage of structural code that T_C executes. Executing the faulty statements is a necessary (but not sufficient) condition to expose software faults. Test coverage measures the test quality in terms of executed code, which approximates the chance to execute a faulty statement (if it exists) [29]. In our experiments, we compute statement and branch coverage, the two most popular code coverage metrics [29].

Simply executing a faulty statement, may not lead to failure. An effective test suite needs test oracles (test assertions) to distinguish correct from incorrect program behaviors. **Mutation analysis** is a well known alternative approach to evaluate the effectiveness of a test suite [20]. Mutation analysis seeds artificial faults in the class under test, producing faulty versions (called mutants). Each of these mutant contains a single-seeded fault. Mutation analysis executes the test suite on each mutant, and counts how many mutants the test suite “kills”, which means that at least one test fails because of the seeded fault. It then computes the mutation score as the percentage of mutants killed by the test suite. Not only does the mutation score evaluate the ability of tests to execute the seeded faults, but also the ability of test oracles to expose such faults.

The values of the test-effort metrics range from zero to one. For example, a line coverage of 0.5, means that a test suite T_C executed half (50 %) of the source code lines in the class under test C .

2.5 Normalization

In this paper we investigate the use of test-quality metrics as normalization factors when measuring the test effort. Current testability approaches measure the test effort by referring only on the size and complexity of the test cases (see Table 2) [12, 13, 46, 51], mostly ignoring the adequacy of the tests (test quality). Focusing only on the size and complexity of test suites without considering their quality may be misleading. A small test suite may reflect the easiness of designing a high-quality test suite, indeed, and thus indicate high code testability, but may also reflect a bad quality test suite, and be completely unrelated to the code testability.

Ignoring the quality of the test suite produces imprecise correlation results, and thus imprecise prediction models if test suites of different quality are considered. For instance, let us consider two classes C_1 and C_2 , and the corresponding test classes T_{C_1} and T_{C_2} . Let us assume that C_1 has 1,000 lines of code (LOC = 1,000), T_{C_1} has ten test cases (T-NOT = 10), and T_{C_1} covers 10% of the source code lines of C_1 (line coverage = 0.1). Let us also assume that C_2 has LOC = 50, T_{C_2} has T-NOT = 30, and line coverage is 90%. Approaches based solely on the size and complexity of the test cases would indicate higher testability (the lower test effort) for C_1 , which requires only 10 test cases for 1,000 lines of code, than C_2 , which requires 30 test cases for 50 lines of code. However, the small size of T_{C_1} comes with very low coverage, which indicates a very low quality of the test suite, while the relatively large size of T_{C_2} comes with very high coverage, which indicates a very high quality of the test suite. The issue is that T_{C_1} and T_{C_2} have a considerably different test quality (10% and 90% line coverage), and thus it is meaningless to compare the number of tests of T_{C_1} and T_{C_2} for studying their correlations with class metrics.

Bruntink and van Deursen’s study is the only work that acknowledges test quality when measuring test effort [13]. The study suggests the importance of test quality, but works around the issue by selecting subjects with test suites that achieves the same code coverage. This approach is feasible when dealing with a small and homogeneous set of subjects, but becomes impractical when dealing with many heterogeneous subjects implemented with different test adequacy criteria (as in our case).

We address this issue by normalizing the test effort with the test quality. We compute both test-effort and test-quality metrics for each test class T_C , and “normalize” the test-effort of T_C with the test-quality of T_C .

Our procedure normalizes the values of each test-effort metrics for all the analyzed systems proportionally to a fixed target test-quality. Our normalization is grounded on the intuition that test effort grows with an increased test quality. We normalize test effort over test quality as:

$$\frac{\text{normalized test-effort value}}{\text{target test-quality value}} = \frac{\text{actual test-effort value}}{\text{actual test-quality value}}$$

The *target test-quality value* is a fixed decimal number between zero (excluded) and one. Intuitively, both code coverage and mutation score is a decimal number within such a range. In our experiment, for simplicity we consider *target test-quality value* to be 1 (but we could have chosen any possible value in the range (0; 1]).

For the example discussed above, with a *target test-quality value* of 1, the normalized value of T-NOT with respect to line coverage is

$\frac{10}{0.10} = 100$ for T_{C_1} and $\frac{30}{0.90} = 33.33$ for T_{C_2} . After the normalization, 1,000 LOCs of C_1 relates with T-NOT = 100 and 50 LOCs of C_2 relates with T-NOT = 33.33, resulting in comparable values.

In the next section, we present our experiments to investigate if such normalization procedure improves the correlation of class metrics with respect to test effort.

3 EXPERIMENTAL RESULTS

This section describes the results of a set of experiments that evaluate our proposed approach for measuring software testability. We addressed two research questions:

RQ1 *What is the correlation between class and test-effort metrics?*

RQ2 *Does the normalization with test-quality metrics increase correlation?*

3.1 Data Collection

We selected 1,186 JAVA open-source projects from GITHUB, the most popular platform for code hosting. We excluded low quality and toy projects, by considering only JAVA projects with at least 50 stars and at least one fork. We queried GITHUB to obtain the list of public repositories that match our criteria. We implemented an automated script that clones the latest version of the master branch for each repository in the list, and selects all the repositories that (i) contain at least one JUNIT test class, which we identify from an import declaration with package `org.junit`, and (ii) use either GRADLE or MAVEN as build automation systems (which we identify from the presence of the file `build.gradle` or `pom.xml`). (iii) build successfully with no failing tests. This is because mutation analysis needs a “green” test suite [20]. We require either GRADLE or MAVEN because we need build automation systems to automatically build the projects, collect and resolve the runtime dependencies and run test cases. Indeed, we need compiled code to compute most class and test-effort metrics (see Table 1 and Table 2), and we need to execute the test cases to compute the test quality metrics (see Table 3). GRADLE and MAVEN are among the most popular build automation systems for JAVA. We found 1,186 GITHUB projects that satisfy these three conditions.

Extracting pairs of class and test class. We automatically explored the content of each of the 1,186 projects to extract the pairs $\langle C, T_C \rangle$, where C is a JAVA class and T_C is the JUNIT test class associated to C . This is in line with the typical usage of JUNIT for unit testing that encodes the tests of a class C in a dedicated class T_C [51]. Following similar work that study the correlation between class and test-effort metrics [12, 13, 24], we identified the pairs $\langle C, T_C \rangle$ by relying on the JUNIT naming conventions for test classes [21, 39]. The conventions require that the name of test class T_C is the name of the associated class C with “Test” or “TestCase” as prefix or suffix [39]. For example, if a class name is `Connector`, the name of its JUNIT test class should be either `ConnectorTest` or `ConnectorTestCase`. Previous studies [12, 13, 24] indicate this as a common practice of JAVA developers. With this approach, we precisely identify the pairs $\langle C, T_C \rangle$ of the class C and the associated test class T_C [12, 13, 24].

Collecting class metrics. For each of the analyzed JAVA classes, we collected the class metrics with a static analyzer that we implemented on top of v2.2 of CKJM-EXTENDED [34] by Jureczko and

Table 4: Descriptive statistics of class metrics

Metric	Mean	SD	Min	Q1	Q2	Q3	Max
LOC	161.68	226.61	8.00	60.00	99.00	177.00	6,758.00
NBI	368.64	917.82	5.00	84.00	183.00	397.00	60,250.00
LOCCOM	149.97	218.04	0.00	52.00	91.00	170.00	6,267.00
NPM	7.98	11.95	0.00	2.00	5.00	9.00	247.00
NSTAM	2.33	6.35	0.00	0.00	1.00	2.00	141.00
NOF	4.50	31.73	0.00	1.00	2.00	4.00	2,083.00
NSTAF	2.21	31.45	0.00	0.00	0.00	1.00	2,083.00
NMC	46.52	111.41	1.00	10.00	22.00	50.00	6,638.00
NMCI	16.10	46.40	0.00	2.00	6.00	16.00	2,533.00
NMCE	30.42	79.16	0.00	5.00	14.00	34.00	4,165.00
WMC	10.50	13.56	1.00	4.00	6.00	12.00	2,560.00
AMC	23.34	43.34	0.11	9.67	16.29	27.00	1,734.50
RFC	29.73	30.36	2.00	12.00	20.00	37.00	520.00
DIT	1.49	0.84	1.00	1.00	1.00	2.00	8.00

Table 5: Descriptive statistics of test-effort metrics

Metric	Mean	SD	Min	Q1	Q2	Q3	Max
T-LOC	112.83	126.18	10.00	50.00	77.00	127.00	3,013.00
T-NOT	5.01	6.98	1.00	1.00	3.00	6.00	191.00
T-NOA	9.26	23.56	0.00	0.00	3.00	9.00	784.00
T-NMC	69.91	116.68	1.00	17.00	36.00	76.00	2,377.00
T-WMC	7.27	8.90	2.00	3.00	5.00	8.00	196.00
T-AMC	34.46	43.69	1.06	15.00	24.25	40.00	1,605.69

Table 6: Descriptive statistics of test-quality metrics

Metric	Mean	SD	Min	Q1	Q2	Q3	Max
Line coverage (L)	0.78	0.26	0.00	0.67	0.88	1.00	1.00
Branch coverage (B)	0.67	0.31	0.00	0.50	0.75	1.00	1.00
Mutation score (M)	0.65	0.32	0.00	0.40	0.71	1.00	1.00

Spinellis [31], currently the most comprehensive open-source tool to compute OO metrics for JAVA. CKJM-EXTENDED computes 18 of the class metrics in Table 1. We implemented additional static analyzers to compute the remaining 11 class metrics. The static analyzers compute the 28 class metrics (Table 1) taking in input the source code of C and the JARs produced with the build automation system, which contain the compiled classes of C 's project and its runtime dependencies. In this way we could compute both the metrics that require the source code of the class C and the ones that require the compiled binary code of the class C .

Collecting test-effort metrics. Table 2 presents the test-effort metrics from T_C that we collected with our static analyzer that already computes LOC, WMC, AMC and NMC. We implemented additional static analyzers for computing the remaining test-effort metrics: T-NOT, and T-NOA. The static analyzer computes the six test-effort metrics by taking in input the source code of T_C and the JARs outputted by the build automation system.

Collecting test-quality metrics. We collected the test-quality metrics (Table 3) relying on the v0.8.2 of JACoCo [42] for the code coverage, and the v1.4.2 of PIT [16] for the mutation score. We built

Metric	Mean	SD	Min	Q1	Q2	Q3	Max
NOC	0.03	0.33	0.00	0.00	0.00	0.00	19.00
MFA	0.21	0.33	0.00	0.00	0.00	0.48	1.00
CBO	6.72	8.17	0.00	2.00	4.00	8.00	156.00
IC	0.33	0.59	0.00	0.00	0.00	1.00	5.00
CBM	0.65	2.04	0.00	0.00	0.00	1.00	48.00
CA	0.71	4.00	0.00	0.00	0.00	1.00	151.00
CE	6.04	7.22	0.00	2.00	4.00	8.00	108.00
LCOM	104.01	755.86	0.00	1.00	6.00	30.00	31,688.00
LCOM3	0.90	0.65	0.00	0.50	0.75	1.10	2.00
CAM	0.43	0.19	0.02	0.29	0.40	0.56	1.00
DAM	0.68	0.44	0.00	0.00	1.00	1.00	1.00
NPRIF	2.74	4.32	0.00	0.00	2.00	3.00	117.00
NPRIM	1.41	3.56	0.00	0.00	0.00	1.00	95.00
NPROM	0.43	1.75	0.00	0.00	0.00	0.00	82.00

an automated script that modifies the `build.gradle` and `pom.xml` build configuration files of each project by adding JACoCo and PIT dependencies. The scripts automatically invokes JACoCo and PIT (via GRADLE v4.10 or MAVEN v3.5.4), which execute each test class T_C individually to compute its line and branch coverage and mutation score. It is worth noting that we used the default configurations for both JACoCo and PIT.

Dataset. We ran the tools and scripts on all the 1,186 project cloned from GITHUB. We aggregated the results by automatically parsing the report files of the static analyzers, JACoCo and PIT. In total, we computed all metrics for 9,861 pairs $\langle C, T_C \rangle$ of class C and associated test class T_C . We counted an average of 8.31 pairs $\langle C, T_C \rangle$ per project. The 9,861 C and T_C classes have a cumulative LOC of 1,594,309 and 1,112,652, respectively. T_C classes have 49,413 test cases overall and 5.01 on average.

Table 4, Table 5 and Table 6 show the descriptive statistics of the values of the class, test-effort, and test-quality metrics in our dataset, respectively. For each metrics, the tables show the average (column "Mean"), standard deviation (column "SD"), minimum value (column "Min"), first quartile (column "Q1"), second quartile (column "Q2"), third quartile (column "Q3") and maximum value (column "Max").

Our dataset contains classes with a wide range of structural and OO design properties (Table 4) and test classes with different size and complexity (Table 5). Table 6 indicates that the 9,861 test classes have a considerable amount of variation of the three test-quality metrics (Standard Deviation (SD) ~ 0.30). This confirms our hypothesis that test quality criteria vary largely among open-source projects, and motivates the need of our normalization adjustment, as discussed in Section 2.5.

The median (Q2) and mean of the test-quality metrics are relatively high, indicating that the projects are well-tested. Line coverage, branch coverage and mutation score have a median of 0.88, 0.75 and 0.71, respectively. This may be related to the selection of popular projects with at least 50 stars and at least one fork. Few test classes result in zero coverage and mutations score (column "Min" of Table 6). In these few cases, the JUNIT naming convention does not correctly pair C with T_C and we excluded these pairs from our analysis.

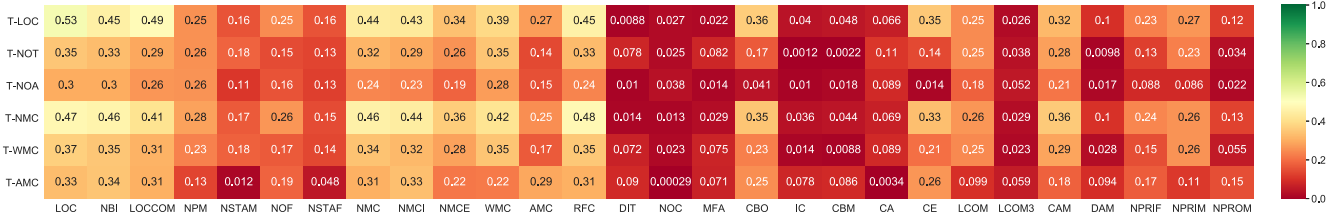


Figure 1: Spearman rank correlation coefficient (absolute values) – without Normalization (N).

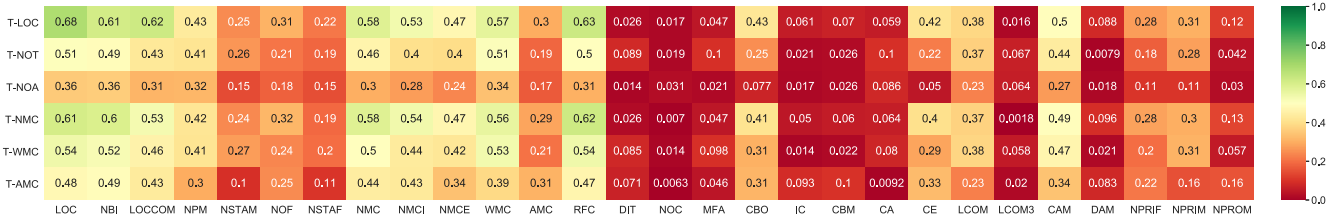


Figure 2: Spearman rank correlation coefficient (absolute values) – test-effort metrics are normalized by Line coverage (L).

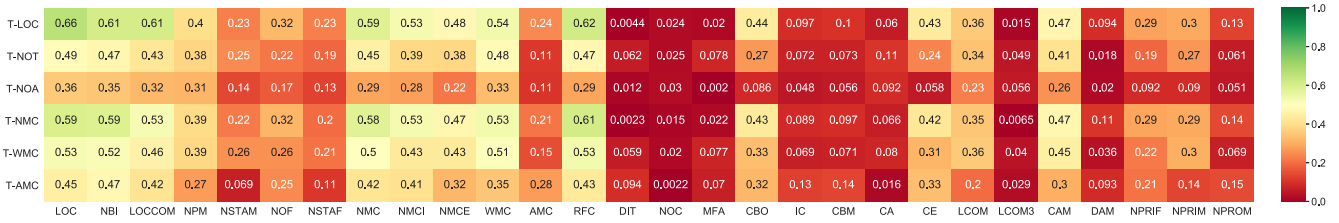


Figure 3: Spearman rank correlation coefficient (absolute values) – test-effort metrics are normalized by Branch coverage (B).

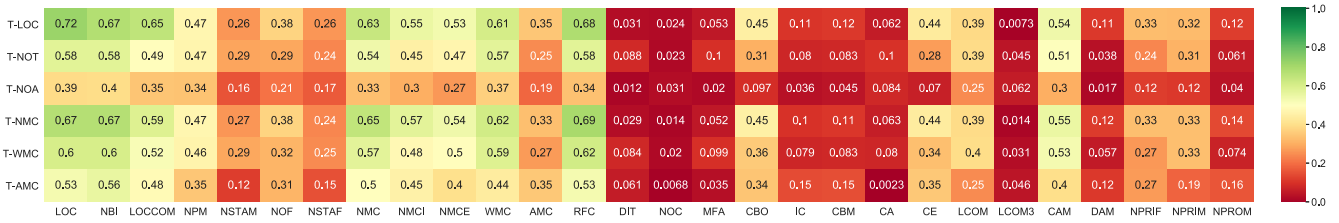


Figure 4: Spearman rank correlation coefficient (absolute values) – test-effort metrics are normalized by Mutation score (M).

3.2 RQ1 – Correlation Study

To answer RQ1, we computed the “Spearman’s correlation coefficient” (ρ) [27] for all 168 (28×6) pair-wise combinations of class and test-effort metrics. *Spearman’s* coefficient is a popular non-parametric measure of correlation used in related studies [13, 51].

We opted for a non-parametric statistical measure because none of the observed metrics (see Table 4, Table 5 and Table 6) follow a “Gaussian” distribution. Thus, parametric measures of correlation, for instance “Pearson”, cannot be used [27]. We checked for normality with the “D’Agostino’s K^2 ” test [19]. Other normality tests, such as “Shapiro-Wilk Test”, are inadequate because each of the distributions has more than 5,000 data points [44]. The K^2 test calculates the *kurtosis* and *skewness* to determine if a data distribution departs

from the normal distribution [19]. For each of the metrics’s value distributions, the normality test leads to $p\text{-value} \leq \alpha$ ($\alpha = 0.05$), and thus we reject the hypothesis that are normal distributions.

Spearman’s coefficient (ρ) quantifies the degree to which two variables are associated with a monotonic function, which is an increasing or decreasing relationship [17]. The coefficient ρ ranges from -1 to $+1$. A positive (negative) ρ means that both variables increase (decrease) together. A ρ close to zero means that the two variables have no correlation.

Figure 1 shows the heatmap of the *Spearman’s* coefficients for each of the 168 pair-wise combinations of class and test-effort metrics. For this research question, we are not interested in distinguishing positive and negative correlations, and thus Fig. 1 reports

absolute values $|\rho|$. The colors range from red (min correlation $|\rho| = 0.0$) to green (max correlation $|\rho| = 1.0$).

We interpret $|\rho|$ as **weak** (≤ 0.3), **moderate** ($0.3 - 0.5$), or **strong** (≥ 0.5), following the widely accepted classification of Cohen [15]. The column “Without Normalization (N)” in Table 7 shows the number of class metrics that have weak, moderate and strong correlations for each test-effort metrics.

The *p*-value (probability value) computed for each moderate and strong correlations is always less than 0.0001, and thus we can reject the null hypothesis that the metrics are uncorrelated. This result confirms those of previous studies [12, 13, 51]. Some class and test-effort metrics have moderate correlations (39 in our dataset).

3.3 RQ2 – Normalization Effect on Correlation

To answer RQ2, we normalized each value of the test-effort metrics with the corresponding value of a test-quality metric, using the formula described in Section 2.5 with 1 as target test-quality value. For example, when considering line coverage, a target test-quality value of 1 means normalizing by 100 % line coverage. Because we considered three test-quality metrics (see Table 3), we obtained three variants of our original dataset. Each variant consider a distinct test-quality metric: Line coverage (**L**), Branch coverage (**B**) and Mutation score (**M**). For each variant, we recomputed the Spearman’s coefficient ($|\rho|$) for the 168 pair-wise combinations of class and (normalized) test-effort metrics.

Figure 2, Fig. 3 and Fig. 4 show the heatmaps of $|\rho|$ after each normalization. Compared with Fig. 1, the values of $|\rho|$ coefficients drastically increase. Table 7 shows the number of weak, moderate and strong correlations after each normalization. Normalizing by Line of coverage (L) decreases the number of weak correlations from 128 to 98 and increases the number of moderate and strong correlations from 39 to 52 and from 1 to 18, respectively. Normalizing by Branch of coverage (B) achieves similar results. Normalizing by Mutations score (M) leads to the best correlation improvement. The number of weak correlations decreases from 128 to 83, while the number of moderate and strong correlations increases from 39 to 51 and from 1 to 34, respectively.

To better quantify the correlation improvement we compared the $|\rho|$ coefficients before and after normalization. We started by removing all combinations of class and test-effort metrics that characterize negligible or not-existing correlations. Removing them is important because their $|\rho|$ values are “statistical fluke” [15], which is a result obtained simply by chance, and not because there is a correlation [15].

Recognize negligible or not-existing correlations is nontrivial because several $|\rho|$ values that are close to zero (weak) before normalization become higher (moderate) after (see Table 7). As such, we removed *all combinations* of class and test-effort metrics that correspond to $|\rho| < 0.1$ *both before and after normalization*, which likely represents statistical flukes [15]. Line coverage normalization has 51 of such combinations, Branch normalization 50 and Mutation normalization 42. The *p*-value for such combinations confirms that the little visible correlation is a *statistical fluke*. Indeed, 91 (63.63 %) of the 143 removed combinations have a *p*-value > 0.0001 , whereas the *p*-values of all retained combinations is < 0.0001 . Therefore, we

can reject the null hypothesis that the retained combinations are uncorrelated.

For each type of normalization, we computed $\Delta |\rho|$ as the difference between the $|\rho|$ values after and before normalization. For example, if we consider the combination LOC and T-NOT, $|\rho|$ before normalization ($|\rho|_N$) is 0.35 and after normalization by mutation score ($|\rho|_M$) is 0.58. Then, $\Delta |\rho| = |\rho|_M - |\rho|_N = 0.58 - 0.35 = 0.23$. The Line normalization retains 117 combinations. Their $\Delta |\rho|$ is 0.09 on average (max 0.17 and min -0.02). Only in four cases the Spearman’s correlation decreases after normalization ($\Delta |\rho| = |\rho|_L - |\rho|_N < 0$). Instead, the Branch normalization retains 118 combinations. Their $\Delta |\rho|$ is 0.08 on average (max 0.18 and min -0.04). The Spearman’s correlation decreases after normalization only in nine cases. The Mutation normalization retains 126 combinations with an average $\Delta |\rho|$ of 0.13 (max 0.26 and min -0.008). The correlation decreases after normalization only for the combination CA and T-NOA: $\Delta |\rho| = |\rho|_M - |\rho|_N = -0.008$. The Mutation normalization leads to the best correlation improvement among the three normalizations, thus sustaining the intuition that the mutation score captures the quality of a test suite better than code coverage [20]. It evaluates both the ability of tests to cover the faulty statements and the ability of test oracles to detect the failure.

3.4 Discussion of the Results

Table 8 shows the Spearman’s coefficients of the 12 class metrics that most correlate with test effort, meaning that they have at least one $|\rho|$ value ≥ 0.4 , either before or after normalization. We chose 0.4 as the threshold because it includes strong correlations and the top half of moderate correlations [15]. Table 8 highlights the highest $|\rho|$ value for each column. All selected class metrics have a positive correlation, except CAM that has a negative one ($\rho < 0$).

These 12 class metrics belong to four OO design properties: size, complexity, coupling, and cohesion (see column “Design Property” of Table 1). More specifically, seven metrics characterize the size of the class (LOC, LOCCOM, NBI, NMC, NMCI, NMCE, NPM), two the complexity (RFC and WMC), two the coupling (CBO and CE), and one the cohesion (CAM). These results confirm some of the findings of related studies [2, 3, 13, 51] and identify new correlations.

We discuss in details the similarities and differences of our findings with respect to Bruntink and van Deursen’s [13] and Toure et al.’s [51] studies, which analyze the highest number of projects and class/test-effort metrics among the work reported in the literature. Referring to the published $|\rho|$ values of these studies, we consider a class metrics to be highly correlated to test effort if (i) the $|\rho|$ value ≥ 0.4 with respect to any test-effort metric; (ii) the correlation results are statistical significant.

Bruntink and van Deursen [13] studied five open-source JAVA programs to compute the Spearman’s correlation between nine Chidamber and Kemerer (C&K) metrics (DIT, NMC, LCOM, LOC, NOC, NOF, NPM, RFC and WMC) and two test-effort metrics (T-LOC and T-NOT). Our study confirms that NMC (FOUT in Bruntink and van Deursen’s study), LOC (LOCC in Bruntink and van Deursen’s study), RFC and WMC highly correlate with both T-LOC and T-NOT. Bruntink and van Deursen also report that NPM is highly correlated with T-NOT for four out of five subjects, and that LCOM and NOF are highly correlated with T-LOC for one subject. Our

Table 7: Counting the Spearman’s coeff. absolute values $|\rho|$ as weak ($|\rho| \leq 0.3$), moderate ($0.3 < |\rho| < 0.5$), or strong ($|\rho| \geq 0.5$)

Test effort	Without Normalization (N)			Normalized by Line coverage (L)			Normalized by Branch coverage (B)			Normalized by Mutation score (M)		
	Weak	Moderate	Strong	Weak	Moderate	Strong	Weak	Moderate	Strong	Weak	Moderate	Strong
T-LOC	17	10	1	15	6	7	14	7	7	11	8	9
T-NOT	23	5	0	17	11	0	17	11	0	15	7	6
T-NOA	28	0	0	21	7	0	23	5	0	19	9	0
T-NMC	17	11	0	14	7	7	14	7	7	11	8	9
T-WMC	21	7	0	15	9	4	14	9	5	13	8	7
T-AMC	22	6	0	16	12	0	17	11	0	14	11	3
Total	128	39	1	98	52	18	99	50	19	83	51	34

Table 8: (Best) class metrics with at least one $|\rho|$ greater than 0.4

Class metric	T-LOC				T-NOT				T-NOA				T-NMC				T-WMC				T-AMC			
	N	L	B	M	N	L	B	M	N	L	B	M	N	L	B	M	N	L	B	M	N	L	B	M
LOC	0.53	0.66	0.66	0.72	0.35	0.50	0.49	0.58	0.30	0.36	0.36	0.39	0.47	0.59	0.59	0.67	0.37	0.53	0.53	0.60	0.33	0.47	0.45	0.53
LOCCOM	0.49	0.60	0.61	0.65	0.29	0.42	0.43	0.49	0.26	0.31	0.32	0.35	0.41	0.51	0.53	0.59	0.31	0.45	0.46	0.52	0.31	0.42	0.42	0.48
NBI	0.45	0.59	0.61	0.67	0.33	0.48	0.47	0.58	0.30	0.36	0.35	0.40	0.46	0.58	0.59	0.67	0.35	0.51	0.52	0.60	0.34	0.48	0.47	0.56
NMC	0.44	0.57	0.59	0.63	0.32	0.45	0.45	0.54	0.24	0.30	0.29	0.33	0.46	0.57	0.58	0.65	0.34	0.49	0.50	0.57	0.31	0.43	0.42	0.50
NMCI	0.43	0.52	0.53	0.55	0.29	0.39	0.39	0.45	0.23	0.28	0.28	0.30	0.44	0.52	0.53	0.57	0.32	0.43	0.43	0.48	0.33	0.42	0.41	0.45
NMCE	0.34	0.46	0.48	0.53	0.26	0.39	0.38	0.47	0.19	0.24	0.22	0.27	0.36	0.46	0.47	0.54	0.28	0.42	0.43	0.50	0.22	0.33	0.32	0.40
NPM	0.25	0.41	0.40	0.47	0.26	0.40	0.38	0.47	0.26	0.32	0.31	0.34	0.28	0.40	0.39	0.47	0.23	0.40	0.39	0.46	0.13	0.29	0.27	0.35
RFC	0.45	0.61	0.62	0.68	0.33	0.49	0.47	0.58	0.24	0.31	0.29	0.34	0.48	0.61	0.61	0.69	0.35	0.53	0.53	0.62	0.31	0.46	0.43	0.53
WMC	0.39	0.55	0.54	0.61	0.35	0.49	0.48	0.57	0.28	0.34	0.33	0.37	0.42	0.54	0.53	0.62	0.35	0.51	0.51	0.59	0.22	0.37	0.35	0.44
CBO	0.36	0.42	0.44	0.45	0.17	0.25	0.27	0.31	0.04	0.08	0.09	0.10	0.35	0.41	0.43	0.45	0.23	0.31	0.33	0.36	0.25	0.31	0.32	0.34
CE	0.35	0.41	0.43	0.44	0.14	0.23	0.24	0.28	0.01	0.05	0.06	0.07	0.33	0.40	0.42	0.44	0.21	0.29	0.31	0.34	0.27	0.33	0.33	0.35
CAM	-0.32	-0.48	-0.47	-0.54	-0.28	-0.43	-0.41	-0.51	-0.21	-0.27	-0.26	-0.30	-0.36	-0.48	-0.47	-0.55	-0.29	-0.45	-0.45	-0.53	-0.18	-0.33	-0.30	-0.40

study does not confirm these results on the large dataset that we used in our experiments.

Toure et al. [51] computed the Spearman’s correlation between seven C&K metrics (CBO, LCOM, WMC, RFC, DIT, NOC, LOC) and three of the test effort metrics considered in our study: T-LOC, T-NOA and T-NMC (TINVOKE in Toure et al.’s study). Our study confirms the high correlation of LOC, RFC and WMC and CBO with test effort, but our experiments do not confirm the LCOM high correlation.

In a nutshell, our experimental results indicate that the OO design properties of size, complexity and coupling largely affect the testability of JAVA classes, confirming the results and conclusions of previous studies [22].

We explain such correlations as follows: *Testability decreases with increasing class size*, as the higher the number of lines of code and methods, the more test cases a developer needs to write. *Testability decreases with increasing complexity*, as effective testing must exercise a number of paths that increases with the complexity of the code. *Testability decreases with increasing coupling*, because classes with high coupling use many components of the system, increasing the execution paths to exercise with the tests [43].

In addition, this paper indicates that also the OO design property of cohesion highly correlates with the test effort: *cohesion is inversely proportional to test effort, thus directly proportional to testability*. Intuitively, low cohesion indicates a low degree of interplay among elements of the same class [45], thus classes with low cohesion need to be exercised with more tests.

3.5 Threats to Validity

A possible threat to external validity is that our results do not generalize to other subjects and OO programming languages. We

mitigated this threat by considering thousands of JAVA projects. The size of our study is several order of magnitude larger than similar studies [12, 13, 51]. Repeat our experiments considering a different OO programming language is an important future work.

A possible threat to internal validity is that there might be errors in our tool or scripts that led to wrong results or metric values. We mitigated this threat by (i) building our static analyzer on top of CKJM-EXTENDED [34], a fully-fledge tool, and (ii) manually validating the correctness of the metric values on a small sample of classes. We release our data and scripts, and welcome external validation [48].

4 RELATED WORK

The definitions of “software testability” [54] can be classified into two groups [22]: (i) “ease of testing”, measured in terms of test-effort metrics (for instance, test size) [12, 13, 46, 51], and (ii) “ease of revealing faults”, measured in terms of test-quality metrics (such as the mutation score and coverage) [1, 18, 30, 33, 57].

In this paper, we comply with the first interpretation, which is the most popular one in literature [22]. However, while we rely on test-effort metrics to measure the ease of testing, we also consider test-quality metrics. Indeed, the key contribution of this paper is to normalize test-effort metrics with test-quality metrics. At the best of our knowledge, this is new to software testability studies.

We discuss the related work on measuring and predicting test effort (which is most closely related to this paper) and test quality, and discuss orthogonal testability work.

Measuring and predicting the test effort. Our work is inspired by the studies of Bruntink and van Deursen on the correlation

between the Chidamber and Kemerer (C&K) and test-effort metrics [12, 13]. Bruntink and van Deursen provide preliminary evidence that C&K and test-effort metrics correlates, and thus the C&K metrics can be used to predict the test effort [12, 13].

Badri and Toure studied the correlation of cohesion metrics (LCOM and LCOM*) with test effort [2] on two projects, concluding that there exists a moderate correlation [2]. Our results do not confirm such a finding. Subsequently, Badri and Toure improved their previous study, increasing the number of software metrics (by adding LOC, CBO, DIT, NOC, WMC, and RFC), and considering three projects instead of two [3]. Badri et al. also investigated the effect of control flow on the test effort [4].

Recently, Toure et al. investigated the use of a metric called “Quality Assurance Indicator” to predict the test effort on eight open-source JAVA projects [51]. Gupta et al. proposed a fuzzy technique to combine values of OO software metrics into a single value called testability index [24]. Singh et al. relied on neural networks to predict testing effort from OO metrics [46].

These studies suffer from three main limitations. First, they involve at most eight software systems. Thus, it is difficult to guarantee that the results generalize to other systems. Notably, a small number of subject systems is common to all testability studies [22]. Among the 182 testability studies that involve the analysis of software systems, Garousi et al. showed 161 (88 %) analyze five or less systems, while the remaining 21 (12 %) studies involve at most 45 systems [22]. Conversely, in our study we analyzed classes from 1,186 projects. Second, they focus on some subsets of the class metrics that we consider in our study. To the best of our knowledge our study is the most comprehensive in terms of the number of class metrics. Third, all of these studies do not consider test-quality metrics for normalization. Lacking a quality assessment of the tests leads to imprecise correlation results and prediction models. Instead, we computed test-quality metrics for each test class and use them to normalize the test-effort metrics.

Measuring and predicting the test quality. Cruz and Eler analyzed four open-source systems, and studied the correlation between Chidamber and Kemerer’s (C&K) metrics and the quality of the tests (coverage and mutation score) [1]. They concluded that the C&K metrics CBO, LCOM, RFC, and WMC have a moderate influence on test quality, and thus a design with low coupling, low complexity, and high cohesion can lead to high coverage and mutation scores.

Khoshgoftaar et al. used neural networks to predict testability based on mutation analysis of source code metrics [33]. Jalbert et al. predicted mutation scores by combining source code metrics with coverage information [30]. Yu et al. proposed a new set of metrics for concurrent programs to predict the mutation score of concurrent tests [57]. Zhang et al. proposed a machine learning approach to predict mutation score from easy-to-access features, for instance, the coverage information, oracle information, and code complexity [58]. Mao et al. extended the approach of Zhang et al. by considering a cross-project setting, more features and more powerful deep learning models [35].

These studies aim to predict test quality, while our work aims to measure and predict the test effort with test-quality metrics (mutation score and coverage) as normalization factors.

Orthogonal work. Design for testability, improvement of testability and fault proneness studies aim at orthogonal goals.

Design for testability aims to measure software testability early in the development process, for example during requirement analysis [7, 32, 53], and design [7–10, 41, 52]. Applying testability analyses early in the development process has the advantage that design refactoring can improve testability before starting the implementation [41].

Improvement of testability aims to refactor programs to increase their testability, for instance, by obtaining a version of the program more amenable to test generation [26].

Basili et al. found that several C&K metrics are associated with fault proneness [6]. Similarly, Gyimothy et al. exploited machine learning methods to predict faults from C&K metrics [25]. Briand et al. explored the relationship between class metrics and the probability of fault detection [11]. Yu et al. examined the relationship between the class metrics and fault proneness [56].

All of these approaches have a different goal with the one of in this paper. An interesting future work is to investigate if our idea of normalizing the test effort with test quality can also help these approaches achieve other testability goals.

5 CONCLUSIONS

This paper proposes a new software testability approach that extends current practice with the novel idea of normalizing test effort with respect to test quality. It also presented the results of an extensive study that involves 9,861 pairs of JAVA classes (with a total of 1,594,309 lines of code) and corresponding JUNIT test cases taken from 1,186 GITHUB projects.

Our results indicate that normalizing test effort with test quality largely increases the correlation between class metric and test effort. An improved correlation between class metric and test effort means a better prediction of test effort.

The normalization procedure that we presented in this paper enables the construction of large-scale prediction models from heterogeneous software systems implemented with different test adequacy criteria. Leveraging our normalization procedure we could train different machine learning models considering different versions of our data-set obtained by normalizing test effort by different *target test-quality values*, such as 70 %, 80 %, 90 % line coverage. Given in input a class, its class metrics values, and a *target test-quality value*, we could predict the test effort using the prediction model corresponding to the *target test-quality value* in input.

In this paper, we introduced the normalization process under the assumption of a proportional growth of the test effort with respect to the test quality. For example, if five test cases (T-NOT = 5) achieve a branch coverage of 50 %, our normalization assumes we need ten test cases (T-NOT = 10) to have a branch coverage of 100 %. One avenue for future work is to study the impact of this assumption on the correlation between class and test-effort metrics.

ACKNOWLEDGMENTS

This work is partially supported by the Swiss SNF project *ASTERIX: Automatic System TEsting of inteRactive software applications* (SNF 200021_178742).

REFERENCES

- [1] Nadia Alshahwan, Mark Harman, Alessandro Marchetto, and Paolo Tonella. 2009. Improving Web Application Testing Using Testability Measures. In *IEEE International Symposium on Web Systems Evolution (WSE)*. 49–58.
- [2] Linda Badri, Mourad Badri, and Fadel Toure. 2011. An Empirical Analysis of Lack of Cohesion Metrics for Predicting Testability of Classes. *International Journal of Software Engineering and Its Applications* 5, 2 (2011), 69–85.
- [3] Mourad Badri and Fadel Toure. 2012. Empirical Analysis of Object-Oriented Design Metrics for Predicting Unit Testing Effort of Classes. *Journal of Software Engineering and Applications* 5, 7 (2012), 513.
- [4] Mourad Badri and Fadel Toure. 2012. Evaluating the Effect of Control Flow on the Unit Testing Effort of Classes: An Empirical Analysis. *Advances in Software Engineering* (2012).
- [5] J. Bansiya and C. G. Davis. 2002. A Hierarchical Model for Object-Oriented Design Quality Assessment. *IEEE Transactions on Software Engineering* 28, 1 (2002), 4–17.
- [6] Victor R. Basili, Lionel C. Briand, and Walcélio L. Melo. 1996. A Validation of Object-Oriented Design Metrics as Quality Indicators. *IEEE Transactions on Software Engineering* 22, 10 (1996), 751–761.
- [7] Benoit Baudry, Yves Le Traon, and Gerson Sunyé. 2002. Testability Analysis of a UML Class Diagram. In *IEEE International Software Metrics Symposium (METRICS)*. 54.
- [8] Benoit Baudry, Yves Le Traon, Gerson Sunyé, and Jean-Marc Jézéquel. 2001. Towards a 'Safe' Use of Design Patterns to Improve OO Software Testability. In *International Symposium on Software Reliability Engineering (ISSRE)*. 324–331.
- [9] Benoit Baudry, Yves Le Traon, Gerson Sunyé, and Jean-Marc Jézéquel. 2003. Measuring and Improving Design Patterns Testability. In *IEEE International Software Metrics Symposium (METRICS)*. 50.
- [10] Robert V. Binder. 1994. Design for Testability in Object-Oriented Systems. *Commun. ACM* 37, 9 (1994), 87–101.
- [11] Lionel C Briand, Jürgen Wüst, John W Daly, and D Victor Porter. 2000. Exploring the Relationships Between Design Measures and Software Quality in Object-Oriented Systems. *Journal of Systems and Software* 51, 3 (2000), 245–273.
- [12] Magiel Bruntink and Arie van Deursen. 2004. Predicting Class Testability Using Object-Oriented Metrics. In *IEEE International Workshop on Source Code Analysis and Manipulation (SCAM)*. 136–145.
- [13] Magiel Bruntink and Arie van Deursen. 2006. An Empirical Study Into Class Testability. *Journal of Systems and Software* 79, 9 (2006), 1219–1232.
- [14] Shyam R. Chidamber and Chris F. Kemerer. 1994. A Metrics Suite for Object Oriented Design. *IEEE Transactions on Software Engineering* 20, 6 (1994), 476–493.
- [15] Jacob Cohen. 2013. *Statistical Power Analysis for the Behavioral Sciences*. Routledge.
- [16] Henry Coles. 2020. *PIT Mutation Testing*. <https://pites.org>
- [17] Gregory W. Corder and Dale I. Foreman. 2011. *Nonparametric Statistics for Non-Statisticians*. John Wiley & Sons, Inc.
- [18] Robinson Crusóe da Cruz and Marcelo Medeiros Eler. 2017. An Empirical Analysis of the Correlation Between CK Metrics, Test Coverage and Mutation Score. In *International Conference on Enterprise Information Systems (ICEIS)*. 341–350.
- [19] Ralph B. D'agostino, Albert Belanger, and Ralph B. D'Agostino Jr. 1990. A Suggestion for Using Powerful and Informative Tests of Normality. *The American Statistician* 44, 4 (1990), 316–321.
- [20] R. A. DeMillo, R. J. Lipton, and F. G. Sayward. 1978. Hints on Test Data Selection: Help for the Practicing Programmer. *Computer* 11, 4 (April 1978), 34–41.
- [21] Mark Fewster and Dorothy Graham. 1999. *Software Test Automation: Effective Use of Test Execution Tools*. ACM Press.
- [22] Vahid Garousi, Michael Felderer, and Feyza Nur Kilicaslan. 2019. A Survey on Software Testability. *Information & Software Technology* 108 (2019), 35–64.
- [23] P. K. Goyal and G. Joshi. 2014. QMOOD Metric Sets to Assess Quality of Java Program. In *International Conference on Issues and Challenges in Intelligent Computing Techniques (ICICT)*. 520–533.
- [24] Vandana Gupta, K.K. Aggarwal, and Y. Singh. 2005. A Fuzzy Approach for Integrated Measure of Object-Oriented Software Testability. *Journal of Computer Science* 1, 2 (2005), 276–282.
- [25] Tibor Gyimothy, Rudolf Ferenc, and Istvan Siket. 2005. Empirical Validation of Object-Oriented Metrics on Open Source Software for Fault Prediction. *IEEE Transactions on Software Engineering* 31, 10 (2005), 897–910.
- [26] Mark Harman, Lin Hu, Rob Hierons, Joachim Wegener, Harmen Sthamer, André Baresel, and Marc Roper. 2004. Testability Transformation. *IEEE Transactions on Software Engineering* 30, 1 (Jan. 2004), 3–16.
- [27] Jan Hauke and Tomasz Kosowski. 2011. Comparison of Values of Pearson's and Spearman's Correlation Coefficients on the Same Sets of Data. *Questiones Geographicae* 30, 2 (2011), 87–93.
- [28] Brian Henderson-Sellers. 1996. *Object-Oriented Metrics: Measures of Complexity*. Prentice-Hall, Inc.
- [29] Joseph R. Horgan, Saul London, and Michael R. Lyu. 1994. Achieving Software Quality with Testing Coverage Measures. *Computer* 27, 9 (Sept. 1994), 60–69.
- [30] Kevin Jalbert and Jeremy S. Bradbury. 2012. Predicting Mutation Score Using Source Code and Test Suite Metrics. In *International Workshop on Realizing Artificial Intelligence Synergies in Software Engineering (RAISE)*. 42–46.
- [31] Marian Jureczko and Diomidis Spinellis. 2010. Using Object-Oriented Design Metrics to Predict Software Defects. Vol. Models and Methodology of System Dependability. Oficyna Wydawnicza Politechniki Wrocławskiej, Wrocław, Poland, 69–81.
- [32] MH Khan and Reena Srivastava. 2015. Flexibility: A Key Factor To Testability. *International Journal of Software Engineering & Applications* 6, 1 (2015), 89.
- [33] Taghi M. Khoshgoftaar, Edward B Allen, and Zhiwei Xu. 2000. Predicting Testability of Program Modules Using a Neural Network. In *IEEE Symposium on Application-Specific Systems and Software Engineering Technology (ASSET)*. 57–62.
- [34] Panagiotis Louridas, Antti Pöyhönen, and Julien Rentropoles. 2020. *CKJM-extended*. http://gromit.iar.pwr.wroc.pl/p_inf/ckjm
- [35] Dongyu Mao, Lingchao Chen, and Lingming Zhang. 2019. An Extensive Study on Cross-Project Predictive Mutation Testing. In *IEEE International Conference on Software Testing, Verification and Validation (ICST)*. 160–171.
- [36] R. Martin. 1995. OO Design Quality Metrics: An Analysis of Dependencies. *ROAD* 2, 3 (1995).
- [37] Michael Mattsson. 1999. Effort Distribution in a Six Year Industrial Application Framework Project. In *IEEE International Conference on Software Maintenance (ICSM)*. 326–333.
- [38] T. J. McCabe. 1976. A Complexity Measure. *IEEE Transactions on Software Engineering* SE-2, 4 (Dec. 1976), 308–320.
- [39] Gerard Meszaros. 2007. *xUnit Test Patterns: Refactoring Test Code*. Pearson Education.
- [40] Joan C Miller and Clifford J Maloney. 1963. Systematic Mistake Analysis of Digital Computer Programs. *Commun. ACM* 6, 2 (1963), 58–63.
- [41] Samar Mouchawrab, Lionel C. Briand, and Yvan Labiche. 2005. A Measurement Framework for Object-Oriented Software Testability. *Information & Software Technology* 47, 15 (2005), 979–997.
- [42] Mountainminds GmbH & Co. KG and Contributors. 2020. *JaCoCo: Java Code Coverage Library*. <https://www.jacoco.org>
- [43] Roger Pressman. 2009. *Software Engineering: A Practitioner's Approach* (seventh ed.). McGraw-Hill, Inc.
- [44] Normadiah Mohd Razali and Yap Bee Wah. 2011. Power Comparisons of Shapiro-Wilk, Kolmogorov-Smirnov, Lilliefors and Anderson-Darling Tests. *Journal of Statistical Modeling and Analytics* 2, 1 (2011), 21–33.
- [45] Robert W. Sebesta. 2012. *Concepts of Programming Languages*. Pearson.
- [46] Yogesh Singh, Arvinder Kaur, and Ruchika Malhotra. 2008. Predicting Testing Effort Using Artificial Neural Network. In *World Congress on Engineering and Computer Science (WCECS)*. 1012–1017.
- [47] Mei-Huei Tang, Ming-Hung Kao, and Mei-Hwa Chen. 1999. An Empirical Study on Object-Oriented Metrics. In *International Symposium on Software Metrics (METRICS)*. 242.
- [48] Valerio Terragni, Pasquale Salza, and Mauro Pezzè. 2020. *Measuring Software Testability Modulo Test Quality - Replication Package*. <https://doi.org/10.5281/zenodo.3740499>
- [49] Fadel Touré, Mourad Badri, and Luc Lamontagne. 2014. A Metrics Suite for JUnit Test Code: A Multiple Case Study on Open Source Software. *Journal of Software Engineering Research and Development* 2, 1 (Dec. 2014), 14.
- [50] Fadel Touré, Mourad Badri, and Luc Lamontagne. 2014. Towards a Unified Metrics Suite for JUnit Test Cases. In *International Conference on Software Engineering and Knowledge Engineering (SEKE)*. 115–120.
- [51] Fadel Toure, Mourad Badri, and Luc Lamontagne. 2018. Predicting Different Levels of the Unit Testing Effort of Classes Using Source Code Metrics: A Multiple Case Study on Open-Source Software. *Innovations in Systems and Software Engineering* 14, 1 (March 2018), 15–46.
- [52] Yves Le Traon, Farid Ouabdesselam, and Chantal Robach. 2000. Analyzing Testability on Data Flow Designs. In *International Symposium on Software Reliability Engineering (ISSRE)*. 162–173.
- [53] Yves Le Traon and Chantal Robach. 1997. Testability Measurements for Data Flow Designs. In *IEEE International Software Metrics Symposium (METRICS)*. 91–98.
- [54] Jeffrey M. Voas and Keith W Miller. 1995. Software Testability: The New Verification. *IEEE Software* 12, 3 (1995), 17–28.
- [55] Michal Young and Mauro Pezze. 2008. *Software Testing and Analysis: Process, Principles and Techniques*. John Wiley & Sons.
- [56] Ping Yu, Tarja Systa, and Hausi Muller. 2002. Predicting Fault-Proneness Using OO Metrics. An Industrial Case Study. In *European Conference on Software Maintenance and Reengineering (ECSMR)*. 99–107.
- [57] T. Yu, W. Wen, X. Han, and J. H. Hayes. 2016. Predicting Testability of Concurrent Programs. In *IEEE Conference on Software Testing, Validation and Verification (ICST)*. 168–179.
- [58] J. Zhang, L. Zhang, M. Harman, D. Hao, Y. Jia, and L. Zhang. 2019. Predictive Mutation Testing. *IEEE Transactions on Software Engineering* 45, 9 (2019), 898–918.