

## Revisiting the Debate: Are Code Metrics Useful for Measuring Maintenance Effort?

Shaiful Chowdhury · Reid Holmes ·  
Andy Zaidman · Rick Kazman

Received: TBD / Accepted: TBD

**Abstract** Evaluating and predicting software maintenance effort using source code metrics is one of the holy grails of software engineering. Unfortunately, previous research has provided contradictory evidence in this regard. The debate is still open: as a community we are not certain about the relationship between code metrics and maintenance impact. In this study we investigate whether source code metrics can indeed establish maintenance effort at the previously unexplored method level granularity. We consider  $\sim 750\text{K}$  Java methods originating from 47 popular open source projects. After considering six popular method level code metrics and using change proneness as a maintenance effort indicator, we demonstrate why past studies contradict one another while examining the same data. We also show that evaluation context is king. Therefore, future research should step away from trying to devise generic maintenance models and should develop models that account for the maintenance indicator being used and the size of the methods being analyzed. Ultimately,

---

Shaiful Chowdhury  
Department of Computer Science  
University of British Columbia, Vancouver, Canada  
E-mail: shaifulc@cs.ubc.ca

Reid Holmes  
Department of Computer Science  
University of British Columbia, Vancouver, Canada  
E-mail: rtholmes@cs.ubc.ca

Andy Zaidman  
Department of Software Technology  
Delft University of Technology, the Netherlands  
E-mail: a.e.zaidman@tudelft.nl

Rick Kazman  
Department of Information Technology Management  
University of Hawaii, Honolulu, Hawaii  
E-mail: kazman@hawaii.edu

we show that future source code metrics can be applied reliably and that these metrics can provide insight into maintenance effort when they are applied in a judiciously context-sensitive manner.

**Keywords** Code Metrics · Maintenance · McCabe · Code Complexity

## 1 Introduction

The cost of software maintenance, which often exceeds the original cost of development [1], has long been a concern for the software industry [2]. This has led to considerable research estimating maintenance effort given the current state of a software project, to support project optimization and risk planning (e.g., [3–9]). External software metrics—such as correctness, and performance—can indicate future maintenance effort, but they are difficult to collect [10] and are often not available in early development phases. In contrast, source code metrics are easy to collect and are available throughout the software development life cycle. Therefore, a holy grail for the developer and the research community has been to predict future maintenance effort from code metrics [10].

A number of code metrics [8, 11–13] have been used to predict maintenance indicators such as defect proneness, change proneness, and test difficulty. However, the true effectiveness of code metrics has been a subject of debate for the past forty years (e.g., [10, 14, 15]). While some studies showed that code metrics were good predictors [16–20], in others the outcome was negative [10, 14, 21]. According to these critics, *other than program size* [10, 22, 23], we do not have a single reliable code metric to estimate software maintenance effort [10, 21, 22]. In fact, size was found to be a good predictor of other code metrics [24], which is frustrating, because if size is the only valid metric, we can not prioritize maintenance activities between two components with similar sizes. Also, no good code metrics except size means that forty years of research [25] on code metrics is potentially useless.

In this paper, we revisit the usefulness of code metrics so that we can inform both the research and developer communities as to whether code metrics are indeed good maintenance predictors, or if they should be abandoned. Also, we reproduce the previous contradictory claims as a means of guiding the research community on how to evaluate future code metrics reliably. For example, while some prior studies accounted for size (usually measured in Source Lines of Code, without comments and blank lines [17]) as a confounding factor for validating a metric, many did not. By using the complete history of ~750K Java methods from 47 popular open source projects, along with six source code metrics and four change proneness based maintenance indicators, we provide encouraging results. Our conclusion is that code metrics can in fact help estimate maintenance effort, such as change proneness, even when the confounding influence of size is eliminated. However, the impact of a code metric varies in different evaluation contexts. For example, *nested block depth* is not as good a predictor for large methods as it is for smaller methods, and

metric performance can vary greatly based on the maintenance indicator used. We support our conclusion by answering the following research questions:

**RQ1:** Is the confounding effect of size a driving factor for the previous contradictory findings on the relationships between code metrics and maintenance effort?

**Contribution 1:** We show that size is indeed a significant factor in previous contradictory claims about the validity of code metrics. With our new method-level benchmark of code metrics and change evolution, we reproduce three major prior observations: 1) Similar to some previous studies (e.g., [16, 26–28]), we first ignore size as a confounding factor, and show that code metrics are good maintenance predictors. 2) By dividing a metric value by size—a common [14, 29, 30], but inaccurate approach [10] for size normalization—we reproduce the claim that code metrics are good maintenance predictors. 3) We then show that the widely adopted size normalization approach fails to neutralize the size influence, and the maintenance impact of code metrics can still be explained by their correlation with size. This reproduces the criticism that without size influence there is no empirical evidence to support the validity of code metrics other than size itself [10, 14].

**RQ2:** Why does the widely used size normalization approach not neutralize the size influence?

**Contribution 2:** Our expectation was similar to many other previous studies: a normalized metric (after dividing by size) should not have any correlation with size. To our surprise, we find that this is not the case. Some normalized metrics are negatively correlated with size while others are positively correlated with size. For example, normalized McCabe values are usually higher when the code size is small, thus producing a negative correlation between maintenance effort and size. For some others, the observation is opposite. We provide an explanation for why this unexpected observation is surprisingly common across all considered metrics.

**RQ3:** Can we apply simple regression analysis for observing the true (size neutralized) maintenance impact of code metrics (proposed in a recent study by Chen *et al.* [31])?

**Contribution 3:** Our conclusion is encouraging. By a combination of bivariate (i.e., size~maintenance) and multivariate (e.g., size + McCabe~ maintenance) regression analysis we show that code metrics are indeed good maintenance predictors, even when their correlation with size is neutralized.

**RQ4:** Does the performance of code metrics vary based on the evaluation context (maintenance indicators and method size), and why?

**Contribution 4:** We show that evaluation context is a significant factor for code metric performance. Some code metrics perform well for small methods, but not for large methods. We show that these metrics lose variability when applied to large methods. Once they reach a threshold they lose predictive power. Other metrics, however, can increase monotonically (e.g., McCabe)

and do not suffer from a lack of variability in the measurements. So their performance is not negatively impacted by code size. Also, a metric’s performance varies greatly based on the maintenance indicator used. A metric can be good for estimating the number of revisions, but not good for estimating the size and the nature of code modifications.

These observations are novel because they clearly show that code metrics are useful as maintenance predictors, while explaining the apparent contradictions from prior studies. With context-based evaluations, we provide new ways to examine the effectiveness of existing and future code metrics, and how they should be used to build more accurate software maintenance models. To aid reproducibility, we provide a public replication package<sup>1</sup> consisting of a data set of  $\sim 750\text{K}$  Java methods with their complete histories and the values for all computed metrics and maintenance indicators over time.

### 1.1 Paper Organization

Section 2 discusses the potential root causes of the previous contradictory claims about code metrics, which helped design the methodology of this paper. Section 3 discusses the methodology. In Section 4, we reproduce the previous contradictory claims about code metrics. We also discuss the inaccuracy of the traditional size normalization approach. In Section 5, we show the true maintenance impact of code metrics by a combination of bivariate and multivariate regression analysis. We also demonstrate why different evaluation contexts should be considered before drawing any conclusion about code metrics. The significance of our findings and threats to validity are presented in Section 6. Section 7 concludes this paper with some potential future studies.

## 2 Related Work & Motivation

First we discuss the McCabe cyclomatic complexity, a metric for measuring the number of linearly independent paths through a component [11]. This metric was proposed in 1976, and has been widely studied and adopted [32,33]. We can divide all the McCabe-related studies into two groups: studies that support its validity (e.g., [4,17,26,34,35]), and studies that do not (e.g., [10,14,21,36]). McCabe is not the only metric that has been debated. With strong empirical evidence, other widely adopted metrics, such as C&K [12], readability [37] have been criticized [10,21]. We identify the following factors that may influence the outcome of a code metric study, and thus support contradictory conclusions.

**User studies are subjective:** Much metrics research relies on user studies to understand the impact of metrics on maintenance indicators [2,21,34,37–43]. One peril related to user studies is that the outcome often depends on human subjects, and can be inconclusive or even contradictory [44]. Also, in the context of code quality, user perception does not necessarily match with the true quality of software [33]. Unsurprisingly, we observe contradictory results for similar maintenance indicators while code metrics were evaluated: both that code metrics are useful [16], and that they are not useful [21]. *In this*

---

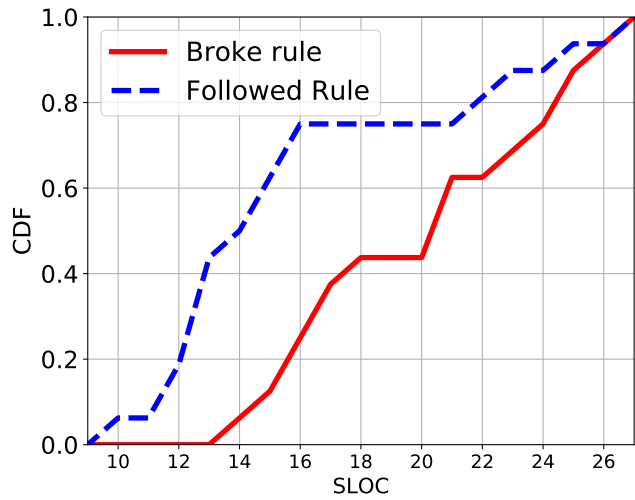
<sup>1</sup> <https://github.com/shaifulcse/codemetrics-with-context-replication>

paper, we therefore focus on objective change measurements as maintenance indicators that we collect from real-world software projects.

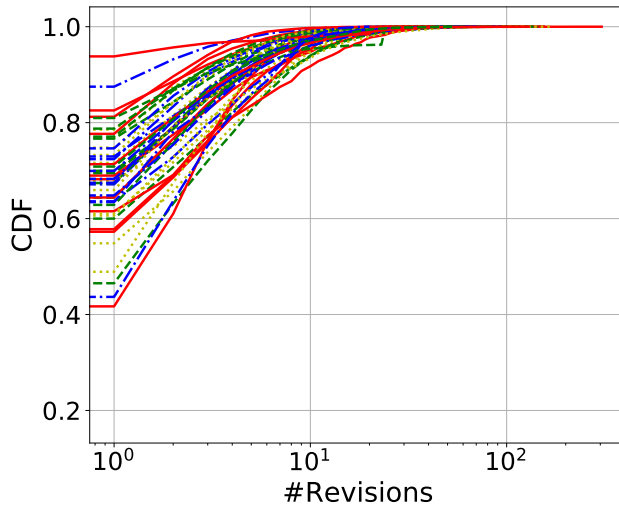
**Size as a confounding factor:** The most frequent criticism invalidating code metrics is that they are highly correlated with size [2, 14, 17, 32, 45]. Therefore, none of the metrics offer any new maintenance information when normalized against size [10, 22]. To claim validity of a metric, we need to show that the metric has predictive power even after its dependency to size is neutralized. Despite this well and long established fact, several studies have ignored it (e.g., [16, 26–28]). In a recent study by Johnson *et al.* [16] developers took less time to read code snippets that followed certain rules (e.g., reduced nesting level) than those that did not follow such rules. The publicly available dataset enabled us to analyze the size distribution of the snippets. Figure 1(a) (cumulative distribution function of source lines of code) shows that the snippets that broke the rule were much larger than those that followed the rules. More lines of code would naturally take more time to read, so perhaps size made the difference in reading time, and not the reduced nesting level.

Some studies, however, have attempted to neutralize size while evaluating code metrics. For example, Spadini *et al.* [18] evaluated the maintenance impact of test smells in three different size categories: small (SLOC < 30), average (30 < SLOC < 60), and large (SLOC > 60). Although this approach should reduce the confounding impact of size to some extent, analyzing all methods with SLOC > 60 (for example) in one group can not eliminate the problem completely. A more common approach is to calculate metric density per lines of code [10, 14, 29, 30]—i.e.,  $100 \times McCabe/Size$ . Unfortunately, Gil *et al.* [10] argued that this approach is inaccurate and questions some of the previous claims of validity for different code metrics. *We argue that a metric is a valid maintenance indicator only when it correlates with maintenance after the confounding factor of size is neutralized, and traditional size normalization approach does not help in making such observation. We need a new approach to evaluate code metrics' effectiveness.*

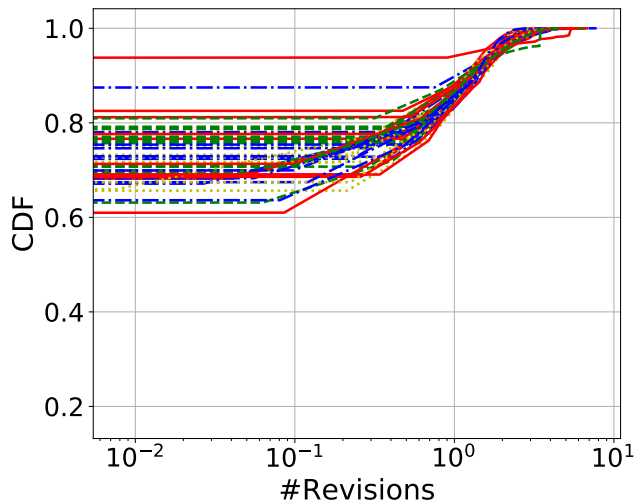
**Aggregated analysis:** Some studies were based on aggregated analyses [10, 18, 46]. That is, they combined all metrics and maintenance indicators from all the studied projects. This is problematic for several reasons. Different external factors—e.g., code review policy [47], developer commit patterns [48] and expertise [49]—cause code to evolve in projects differently. Figure 1(b) shows the distributions of revisions for all the methods in each of our 47 projects (described later); each line corresponds to one project. Evidently, these projects do not exhibit similar revision behavior. Combining them together may lead us to inaccurate conclusions. Figure 1(c) shows that the difference in distributions reduces after applying recommended log-normal transformation [15], but the differences do not completely disappear. Even for code metrics, the distribution in their measurements greatly vary based on a project's domain, programming language, and life span [50]. Also, some projects are much bigger than others. This means that in aggregated analyses, results can be unduly influenced by few big projects.



(a) Confound factor



(b) Distribution of method revisions



(c) Distribution after log-normal transformation

Fig. 1: Figure (a) shows that size was not normalized in the study by Johnson *et al.* [16]. Code snippets that broke the rules are much larger than the code snippets that followed the rules, leading to inaccurate comparisons. Figure (b) shows that aggregated analysis is inaccurate because different projects exhibit different revision behaviors. Each line represents the revision distribution for a given project, and these lines are very different from each other. Figure (c) shows that these different revision patterns are not neutralized, even after applying a log-normal transformation, as suggested by Gil *et al.* [15].

These problems of aggregated analysis can be avoided by analyzing each project individually [2–4, 28]. Individual project analysis, however, has been criticized for selection and publication bias [10, 51]. As we also show in this paper, there are always outlier projects that exhibit unique behavior, which might seem normal if too few projects are studied. *The argument is thus to analyze each project separately, while studying a reasonably large number of projects with a systematic unbiased selection process.*

**Granularity:** Software maintenance studies have been conducted at different granularities that can influence observations [17]; these include the system level [2, 29], class/file level [52], snippet level [39, 42], and even `git diff` level [53]. Understanding maintenance at the method level granularity from real software evolution data is difficult [54–56]; it is harder to reliably generate method level histories than file level histories. Despite the difficulty, method level is the most desirable granularity [46, 57], because class/file level granularity is often too coarse-grained for practical use [46, 58, 59]. *We also argue that if we can estimate maintenance at method level granularity, we can extend this understanding to coarser levels of granularity—a class is generally a collection of methods.*

**Maintenance Indicators:** We consider software maintenance as a *construct* [60], which is difficult to measure, but easier to estimate it through some reflective indicators. Different studies have focused on different indicators: human effort to read and understand code [16, 61], localizing bugs [9, 46, 57, 62, 63], change proneness [10], or developer activities [3]. In this paper, we focus on four change proneness indicators of Java methods (justified later). We show that metric performance to understand maintenance can vary significantly based on the indicators used.

### 3 Methodology

This section describes our process for: i) selecting projects, ii) choosing code metrics and maintenance indicators, iii) collecting method-level history for analysis, iv) age normalization for methods with different ages, and v) selecting statistical approaches for analysis.

#### 3.1 Project Selection

To reduce inaccuracies that may stem from aggregated analysis, we opted to analyze individual projects. To neutralize selection bias, we took the union of all GitHub Java projects used in four different software evolution studies [10, 18, 52, 56], totalling 47 projects—mixing projects from different programming languages can significantly impact the outcome of code metrics studies [50]. As we show later, this project set is able to highlight code metric behaviors that are generic (true for most projects) and which behaviors are rare. Table 1 describes the dataset. The table suggests that only a small number of methods (e.g., 95<sup>th</sup> percentile of revisions) undergoes a large number of revisions. That means we can significantly reduce the search space for maintenance optimization by identifying the top 5% high-churn methods. This paper investigates, if code metrics are indeed helpful for such identification.

Table 1: Description of the dataset used in this paper. In total, 744,197 Java methods were collected from 47 GitHub Java projects. The small average, rev (avg), and median number of revisions compared to the large 95<sup>th</sup> percentile and maximum revisions suggest that most maintenance activities occur in small areas of code.

Repository	# methods	# rev (avg)	# rev (med)	# rev (max)	# rev (95 <sup>th</sup> percentile)
hadoop	70,524	1.8	1.0	67.0	6.0
elasticsearch	63,253	3.7	2.0	133.0	12.0
flink	38,482	1.9	1.0	96.0	8.0
lucene-solr	37,600	1.6	1.0	145.0	6.0
hbase	37,414	3.2	2.0	109.0	11.0
docx4j	36,543	2.2	2.0	49.0	4.0
intellij-community	36,387	3.6	2.0	158.0	13.0
weka	35,824	1.7	1.0	86.0	5.0
hazelcast	35,609	2.7	1.0	109.0	10.0
spring-framework	27,093	2.4	1.0	60.0	8.0
hibernate-orm	25,290	2.5	2.0	70.0	7.0
eclipseJdt	22,230	3.0	1.0	133.0	12.0
jclouds	21,478	1.7	1.0	65.0	6.0
guava	21,195	1.1	0.0	45.0	4.0
sonarqube	20,907	3.1	2.0	305.0	10.0
wildfly	19,980	2.2	1.0	94.0	9.0
netty	17,089	2.1	1.0	75.0	9.0
cassandra	16,271	1.6	0.0	178.0	6.0
argouml	12,775	3.4	2.0	156.0	10.0
jetty	10,743	2.2	1.0	93.0	8.0
voldemort	10,654	1.8	0.0	65.0	9.0
spring-boot	10,495	2.6	2.0	59.0	9.0
wicket	10,260	4.9	3.0	63.0	14.0
ant	9,837	2.0	1.0	73.0	8.0
mongo-java-driver	9,735	3.4	2.0	57.0	13.0
jgit	9,615	1.4	1.0	44.0	6.0
pmd	9,037	3.2	2.0	91.0	10.0
xerces2-j	8,335	1.4	0.0	65.0	6.0
RxJava	8,332	3.7	3.0	22.0	10.0
javaparser	6,324	3.4	1.0	84.0	15.0
openmrs-core	6,108	2.2	1.0	176.0	7.0
hibernate-search	5,391	3.2	2.0	61.0	11.0
titan	4,654	2.2	1.0	42.0	9.0
checkstyle	3,347	3.8	2.0	72.0	13.0
commons-lang	2,964	3.7	3.0	34.0	9.0
lombok	2,852	2.5	1.0	61.0	10.0
atmosphere	2,681	2.4	0.0	87.0	12.0
jna	2,653	2.3	1.0	38.0	8.0
Essentials	2,412	3.1	1.0	46.0	15.0
junit5	2,181	2.7	1.0	59.0	11.0
okhttp	1,970	4.7	3.0	54.0	16.0
hector	1,963	1.7	1.0	43.0	8.0
mockito	1,526	4.2	3.0	62.0	13.0
cucumber-jvm	1,222	2.7	1.0	36.0	9.0
commons-io	1,149	3.1	3.0	24.0	8.0
vraptor4	934	1.7	1.0	24.0	6.0
junit4	879	3.1	2.0	70.0	11.1



We also note that the number of methods is significantly different across the projects. If we were to adopt aggregated analysis, results for small projects would be unnoticeable. The set of projects is clearly diverse. For example, even for the subset that was used in [10], the number of developers ranges between 16 and 197, and development duration varies from one to 13 years.

### 3.2 Code Metric Selection

In contrast to method-level granularity, many of the popular code metrics, such as C&K and depth of inheritance, work only at class or higher level granularities. Also, the objective of this paper is not to show which code metric is the best for estimating maintenance effort, because there are many of them [21]. Instead, we focus on code metric validity: are they useful at all, and if so, are the underlying evaluation contexts important? Therefore, we focus only on six widely adopted and widely studied metrics that are applicable at method-level granularity. We show that these six metrics were sufficient to reproduce previous contradictory claims and to examine how code metrics should be reliably evaluated.

**McCabe:** The McCabe algorithm for measuring cyclomatic complexity is simply:  $1 + \#predicates$  [11]. There are, however, two forms: one counts logical `&&` and `||`, and the other ignores them. We only consider the latter form because considering them does not make any meaningful difference in McCabe’s validity as a code metric [17].

**McClure:** A criticism of McCabe is that it does not consider the number of control variables in a predicate. If the outcome of a predicate depends on multiple control variables, it should be considered more complex than the one with a single control variable [2]. McClure differs in this regard [8]: it measures the sum of the total number of comparisons (thus includes `&&` and `||`) and the number of control variables in a component.

**Nested Block Depth:** Neither McCabe, nor McClure, considers nesting depth. To both of these metrics, two methods each with two loops (for example) are equally complex, even if one of them has nested loops and the other does not. Measuring Maximum Nested Block Depth (referred to as NBD) is a common solution [16, 64].

**Proxy Indentation:** Hindle *et al.* [53] argued that a metric like McCabe is hard to calculate because one needs a language-specific parser. They found that it is similarly useful to use the level of indentation in a code component. Counting the raw number of leading spaces in each line is equally good as counting the number of logical spaces. Instead of calculating the max, sum, mean, or median, the authors found that standard deviation of those counts (referred to as IndentSTD) works as the best proxy for McCabe-like complexity.

**FanOut:** The aforementioned metrics, to some extent, measure similar complexities—mainly the number of conditional branches. Therefore, we add FanOut (total number of method calls made by a given method) to our list.

This metric provides an indication of how a particular method is dependent on other methods (coupling). Mo *et al.* observed that highly coupled systems are usually less maintainable [65]. We also wanted to use FanIn or unique FanOut, but these two require a symbol solver that preprocesses a complete repository for each change commit a method has. It would be extremely time demanding for the  $\sim 750\text{K}$  methods that we consider.

**Readability:** We also include a composite metric that combines different code metrics to produce a single indirect maintenance index. For this we adopt the widely used Readability metric by Buse *et al.* [37] which ranges from 0 (least readable) to 1 (completely readable code).

### 3.3 Maintenance Indicator Selection

Modeling maintenance effort is a difficult problem, because there are many different effort indicators that should be considered for building a comprehensive effort prediction model. A subset of these indicators include human effort to read and understand code [1, 16, 21, 37], difficulty to modify a code [65], bug proneness [59, 66], and change proneness [10, 28, 58, 67]. The objective of this paper is not to build an effort prediction model, but to answer if code metrics are at all useful for understanding maintenance effort, and how to evaluate these metrics reliably. In that vein, we focus on change proneness, as it is measurable without conducting user studies, reducing threats related to such studies. Also, the community unanimously agrees about the utility of change proneness as one of the most applicable maintenance effort indicators [10, 28, 52, 54, 58, 67–69]. While we considered bug proneness, we discarded this indicator to reduce threats to construct validity. From our dataset, bug proneness can be measured by capturing keywords from commit messages, such as *error*, *bug*, and *fixes* [68, 70]. Unfortunately, this approach has been criticized for low precision/recall [10, 18], which is further complicated due to tangled changes. Developers often commit unrelated changes, which incorrectly labels bug-free code as *buggy* [48]. Additionally, change proneness is often highly correlated to bug proneness [46, 66, 71–73]. Therefore, if a code metric is a good predictor of change proneness, it is likely to be a good predictor of bug proneness as well. Ultimately we selected the following four change proneness indicators.

**#Revisions:** Number of revisions of a component is considered as an indication of maintenance effort by many [3, 20, 64, 74]. The consensus is that a well designed less complex component should not need many revisions.

**Diff size:** Number of revisions does not disclose how large a change is. If two components are revised the same number of times, their maintenance effort is not necessarily the same. Also, the number of revisions can be influenced by developers' commit habit or culture [47]. Therefore, some consider `git diff` size a more accurate maintenance indicator [3, 75].

**Additions only:** Adding new lines is perhaps more difficult than deleting lines. This threat can be reduced by considering only the number of new lines added [3].

**Edit Distance:** Lines of changed code, as a metric, is affected by noise such as coding style; it does not distinguish modifications between large and small lines. Also, a simple automatic rename method refactoring may modify a large number of lines. Therefore, Levenshtein edit distance [76] is considered as a better maintenance indicator than number of lines (added and/or deleted) [5, 21, 75]. Levenshtein edit distance measures the number of characters *added* + *deleted* + *updated* for converting one source code version into another.

### 3.4 Data Collection and Representation

We require a method's complete change history: how many times the method was changed, when the changes happened, and what was changed? There are only few tools that support history tracing at method level granularity: Historage [77], FinerGit [55], and CodeShovel [56, 78]. Historage and FinerGit work similar to `Git`'s file tracking mechanism by converting each Java method to a file. However, we find that this approach does not scale well to larger projects. In contrast, CodeShovel tracks a method (even if the method's signature is changed) using string similarity and without any project preprocessing. Unlike the other tools, CodeShovel's accuracy was evaluated on both open source and closed source industry projects, with 99% precision and 90% recall.

After collecting the complete history of 744,197 Java methods from 47 selected projects, we collected the evolution of their code metrics (e.g., SLOC, McCabe), and change metrics (e.g., edit distance). To the best of our knowledge, there is no existing tool that provides measurements in this form, so we have implemented our own tool. We verified its correctness by randomly selecting and validating 200 Java methods. In addition, the accuracy of the tool was tested by an independent code metric researcher. A method, across its evolution history, can have different values for the same code metric (e.g., initially the McCabe was 5, but then it changed to 3, and then to 5 again). For a single method, we thus summed all the maintenance indicator values (e.g., sum of all edit distances) that a method had for each unique code metric value. For a given method, for example, if edit distance 10, 20, and 30 correspond to McCabe values 5, 3, and 5 respectively, McCabe value 5 is blamed for edit distance 40 (10+30), and McCabe 3 is blamed for edit distance 20. This is how we mapped code metrics value with different maintenance indicators to study the relationships between them.

### 3.5 Age Normalization

It is inaccurate to compare the change history of two differently-aged methods. An older method is more likely to have more revisions than a newer method [45, 79]. For the rest of the analysis, we consider methods that are at least two years old: reasonably enough time to undergo their initial changes. However, this approach does not completely neutralize the time effect; for instance, we should not compare a two year with a year ten method. We neutralize this by considering changes that happen only within the first two years of these filtered samples. This is like time traveling to each of the methods

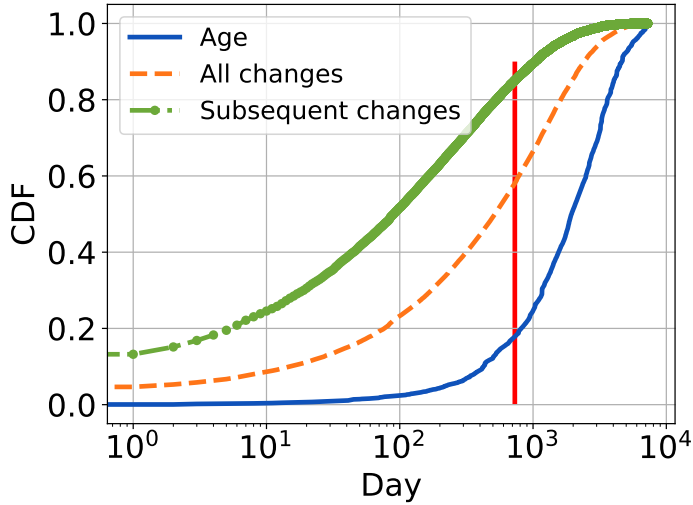


Fig. 2: Two years is a good threshold for age normalization. We only lose  $\sim 20\%$  of the methods, and yet retain  $\sim 60\%$  of the revisions that happened within our whole dataset.

change history when they were two years old. But why two years? Figure 2 (cumulative distribution functions with day) shows that more than 80% of our methods (total 611,543 methods) are older than two years (Age). Among all the revisions in whole dataset (All changes in the graph),  $\sim 60\%$  of them happened within the first two years. If we consider the interval time of subsequent revisions, around 86% of changes happened within the first two years. If we increase the age threshold value, we lose more methods. If we decrease it, we lose more change history, so it is a trade-off. Note that if we set the threshold to one or three years, the major conclusions of this paper remain the same.

### 3.6 Correlation and Statistical Significance

To apply Pearson’s formula for calculating correlation coefficients between code metrics and maintenance indicators we need to establish that each metric, for each change proneness indicator, for each project is normally distributed. After applying the Anderson-Darling normality test [80] for some of the randomly selected projects, we found that they are not normally distributed. Therefore, we opted to use Kendall’s  $\tau$  correlation coefficient. Unlike Pearson’s correlation coefficient, Kendall’s  $\tau$  does not assume any distribution of the data (non-parametric), and is less affected by outliers, which the community has chosen to use for these kinds of analysis [10, 81, 82]. Unless otherwise stated, all results in this paper are statistically significant (p-value  $< 0.05$ ). When necessary, we also use the Wilcoxon rank-sum test to test if the performance distributions of the code metrics are statistically different, and if so, we report how large the differences are (Cliff’s Delta effect size). Similar

to Kendall’s  $\tau$ , these two tests are non-parametric and do not assume any distribution of the data [83, 84].

As we consider each project separately, we present the results as distributions. Therefore, we use the Cumulative Distribution Function (CDF) for the visual representation of our results. We considered using XY-plots, but CDF better conveys our findings. As CDF is a monotonic function, comparing multiple lines (because of multiple code metrics) is easier than XY-plots’ zigzag-patterns.

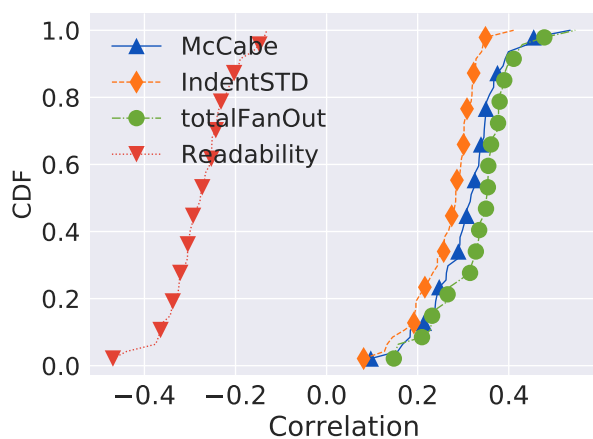
#### 4 Results: Looking into the Past

In this section, we reproduce previous claims about the relationship between code metrics and software maintenance (RQ1). We show that the debate about source code metric effectiveness stems from improperly considering, or normalizing for, size as a confounding factor. We show that the most commonly used normalization approach fails to neutralize the size effect in practice. We then explain why size normalization is difficult and remains an open research problem (RQ2).

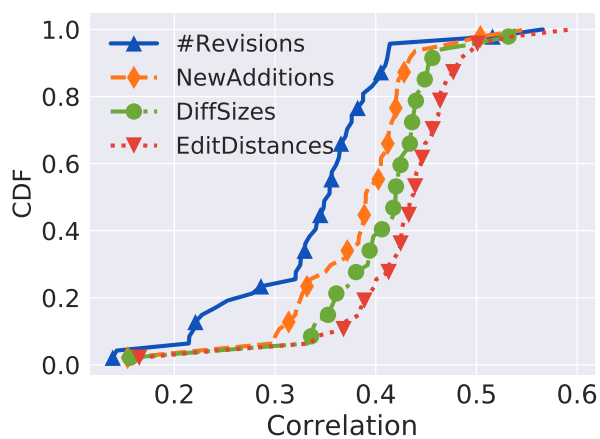
##### 4.1 (RQ1) Metrics are (not) Useful

Figure 3 (a) shows the cumulative distribution functions (CDF) of the correlation coefficients between the selected code metrics and number of revisions (each line represents a particular metric and shows the distribution across all the 47 projects). Results are similar for McCabe, McClure, and NBD, so we show only McCabe to maintain graph readability. Evidently, all six code metrics are correlated with maintenance measures, which means that they are potentially good maintenance predictors. This approach aligns with a group of prior studies (e.g., [16, 26, 27]) that did not attempt to control for size as a confounding factor.

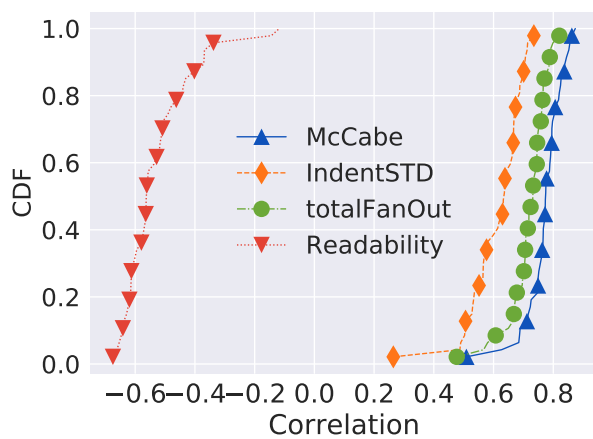
Correlating a metric with maintenance alone does not make a metric valid or useful [10, 14]. The arguments supporting this are: (1) size is a great predictor of maintenance, and (2) many code metrics are highly correlated with size. So a metric’s correlation with maintenance could simply be due to its correlation with size. Figure 3(b) shows the correlation between SLOC and the four maintenance indicators for all 47 projects (supporting argument 1). Figure 3(c) shows that all the metrics are correlated with size (supporting argument 2). We observe that the direction and strength of the correlation between a code metric and maintenance is similar to the metric’s correlation with size. For example, Readability is negatively correlated with both size and maintenance (larger size means less readable and thus less maintainable). For IndentSTD the correlation is somewhat lower (compared to McCabe and FanOut) with size and thus lower with maintenance. These observations align well with the criticism that when the influence of size is considered, we do not have any empirical evidence to support the usefulness of code metrics [14]. From this we conclude that without size normalization we do not know the true effectiveness of code metrics.



(a) Code metrics and #Revisions.



(b) SLOC and change indicators.



(c) Code metrics and SLOC.

Fig. 3: Figure (a) shows that all the code metrics are significantly correlated with number of revisions in each project (observations are similar for other maintenance indicators). Figure (b) shows that SLOC is positively correlated with all the maintenance indicators. Figure (c) shows that code metrics are correlated with SLOC. For graph readability, the number of marks in each line is fewer than the actual number of data points.

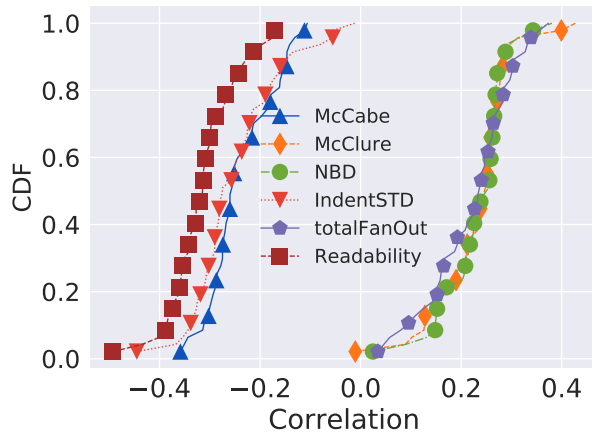
The most common approach for size normalization takes the density of a metric and divides its measure in a component by the size of the component [10, 14, 29, 30]. For example,  $100 \times McCabe/Size$  gives the McCabe value per 100 lines of code, so we should have a normalized McCabe measure completely independent of size. The hypothesis is that, if we still see correlation between a metric and maintenance, we can argue for the validity of the metric. Figure 4(a) shows the distributions (for 47 projects) of correlation coefficients for all the normalized code metrics with the number of revisions (results are similar for other maintenance indicators). Evidently, all the metrics are still correlated with maintenance. This supports the assumption of code metrics validity after the size influence is neutralized.

To our surprise, we find that this commonly practiced size normalization approach is inaccurate. Figure 4(b) shows the correlation distributions between the normalized metrics and size for all 47 projects. Although we were expecting the correlation to be close to zero, this is not the case. We also note that the direction of the correlation between size and a metric still dictates the direction of the correlation between maintenance and that metric. For example, McCabe is negatively correlated with both size and number of revisions. We later found that the same observation was made by Gil *et al.* [10], although their granularity level was different (file-level instead of method-level). They concluded that *size is the only valid code metric* because maintenance impact of other code metrics can directly be explained by their correlations to size. Despite our similar observation, we see hope if we carefully examine Figure 4(a) and 4(b). For example, the correlations with revisions are similar for FanOut, NBD, and McClure, but not as similar to size. For Gil *et al.*, size explains everything (file-level), but for us it does not (method-level).

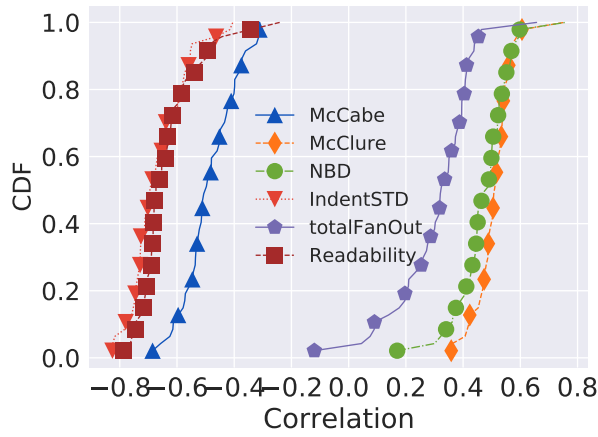
This difference of observations can be explained by the findings of Landman *et al.* [17], who have studied the correlation between McCabe and size at different granularity levels. In their study, the strong correlation between McCabe and size is true only for large code units, but dwindles significantly at the method level granularity. The authors, however, did not examine McCabe’s impact on maintenance. In this paper, with the help of bivariate and multivariate regression analysis, we show that code metrics are indeed good maintenance predictors, even when their relations with size are neutralized. But before examining this, we first explain why the widely adopted and believed size normalization fails (RQ2).

#### 4.2 (RQ2) Size Normalization is Sensitive to Size

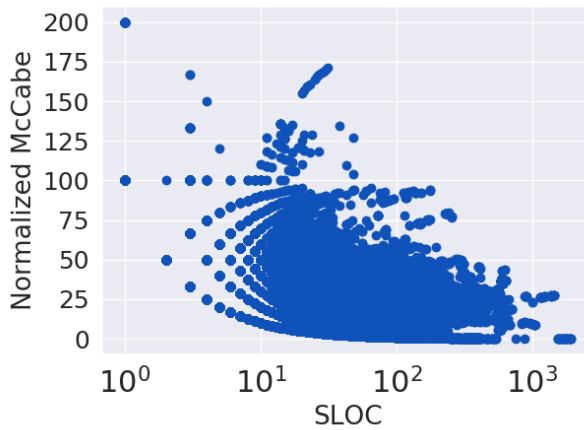
To see why size normalization does not work, we take a deeper look into McCabe complexity ( $1 + \#predicates$ ). Why is normalized McCabe negatively correlated with size (and thus with maintenance)? Of course, a negative correlation here indicates that the lower the size the higher the normalized value (i.e., density per 100 lines of code). Interestingly, we find that the ‘1’ in the McCabe formula ( $1 + \#predicates$ ) is a major issue. Consider a simple 3 line Java method, which just returns the sum of two numbers. The McCabe is



(a) Normalized metrics and #Revisions.



(b) Normalized metrics and SLOC.



(c) Normalized McCabe against SLOC.

Fig. 4: Figure (a) shows that all the normalized code metrics are still significantly correlated with number of revisions in each project. However, Figure (b) suggests that even the normalized code metrics are correlated with SLOC. Figure (c) demonstrates why the normalized McCabe is negatively correlated with SLOC (and thus with change).



```

static int normalizeIndex(final byte[] ar,
final int index) { return index >= 0 ?
index : index + array.length; }

```

Fig. 5: A sample method from Elasticsearch.

already 1, and the normalized value is 0.33 (1/3). The effect of ‘1’, however, diminishes as method size grows. Figure 4(c) depicts normalized McCabe score against size for all the 611,543 methods. Clearly, the normalized McCabe score decreases with the increase in size because of the high density in small methods. We find that if we eliminate the ‘1’ from the formula, the size influence in normalized McCabe is reduced, but does not go away completely. The graph shows that the normalized score for small methods can even exceed 100. As an example, consider the method from the *Elasticsearch* project, shown in Figure 5. This method is written as a one line method (size = 1), and because of the conditional operator, the plain McCabe is 2 (1+#predicates). The normalized McCabe is thus 200 (100×(2/1)). Note that this data point with 200 normalized McCabe score represents 56 methods, not just one.

As we show later, most methods in our dataset are small (SLOC  $\leq 22$ ), and therefore the overall correlation is significantly impacted by the high McCabe density of this large number of small methods. The problem is, Figure 4(b) shows that different code metrics suffer differently from this size normalization approach. Size normalization thus remains an open research problem: we need to develop an approach that not only eliminates the influence of size, but also does not normalize in a way that hides the effectiveness of code metrics.

**Summary:** Previous approaches that supports the usefulness of code metrics either did not consider the size influence or normalized it inaccurately. A new way is required to evaluate the true effectiveness of code metrics by completely eliminating the effect of size.

## 5 Results: Evaluating metrics with regression and contexts

The problems of performing a size-decoupled metric evaluation with a traditional size normalization approach led us to a study by Chen *et al.* [31]. The authors investigated why different mutation testing studies claimed differently [81,85–87] about the relationship between test suite size, test adequacy criteria (e.g., coverage), and test effectiveness (fault detection). Although the context of their study is different than ours, the outcome is similar: two different studies control for test suite size, while evaluating the relationship between test adequacy criteria and test effectiveness, producing two different conclusions. For highly correlated variables (code metrics and size in our case), the authors suggested that regression analysis can be useful. Encouraged by their hypothesis, we designed our approach as follows.

1. With a bivariate regression analysis between size and a maintenance indicator (e.g., revisions), we calculate the goodness of fit score of the regression model.

2. In the same model, we then add one of the code metrics (e.g., McCabe) as the second independent variable and asked whether this multivariate regression model ( $\text{size} + \text{McCabe} \sim \# \text{revisions}$ ) improves the goodness of fit score with statistical significance ( $p\text{-value} < 0.05$  for the coefficient of McCabe)?
3. We take the difference between the two fitness scores and convert it to a percent improvement to show the distribution across the 47 different projects.
4. We repeat steps 1 to 3 for all the maintenance indicators and code metrics.

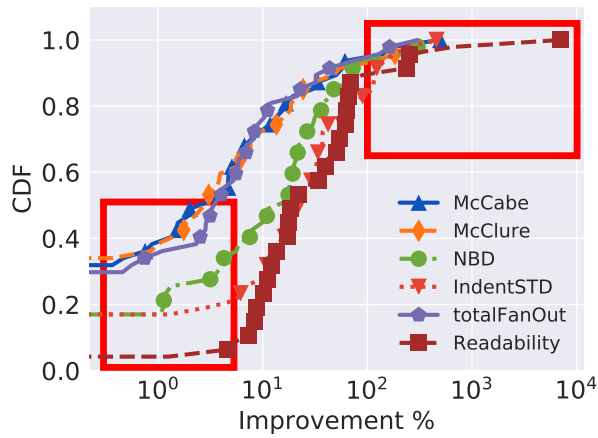
We argue that this approach shows the true maintenance impact of code metrics because it correctly eliminates the size influence.

### 5.1 (RQ3) Regression Analysis for Code Metric Evaluation

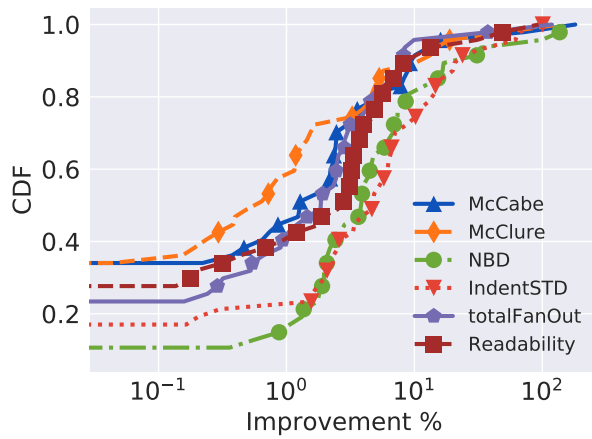
Figure 6 shows the CDFs of percent of improvement in fitness scores with the multivariate regression models for three of the maintenance indicators. We excluded the result for the `diff size` indicator, because the observation is the same—all the metrics improve the fitness scores for all of the indicators. Figure 6(a), for example, shows that the Readability metric improves the fitness accuracy by at least 10% for more than 80% of the projects (and at least 100% for 20% of the projects). This clearly refutes the claim that code metrics are not useful after size influence is neutralized [10, 22]. The performance of these metrics, however, are not the same across all maintenance indicators. Readability is the best metric for estimating the number of revisions. For estimating change size (e.g., edit distance), however, NBD and IndentSTD are the best estimators. Also, results in Figure 6 are dominated by methods that are small in size, because most of our methods are small. We need to evaluate if code metrics perform differently when evaluated for large methods only, and if so, what factors influence their performance. These are the questions we investigate in RQ4.

### 5.2 (RQ4) Evaluating Code Metrics with Contexts

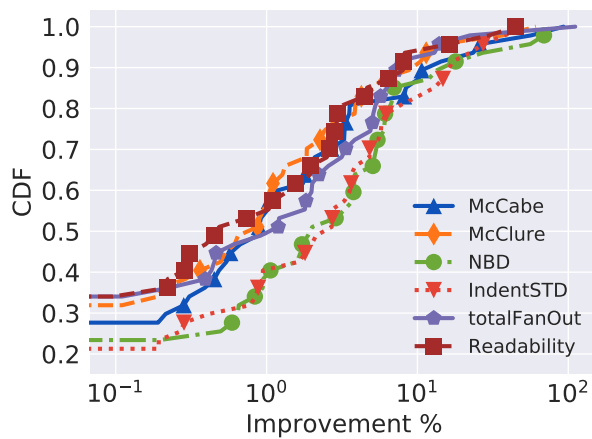
To investigate whether the performance of a code metric depends on the method size, we need to first define a threshold for separating large methods from small ones. Instead of defining such a threshold from intuition or from expert opinion, we followed the 6-step systematic approach proposed by Alves *et al.* [88]. The main objective of the approach is to find critical values for identifying low risk (small size), medium risk (medium size), high risk (large size), and very high risk (very large size) code components in terms of maintenance from a given set of projects. These critical values are robust, i.e., they are not impacted by outlier projects or methods. We refer to [88] for more detail. The first 5 steps of Alves *et al.*'s approach deliver 3 critical values that are derived from Figure 7. The first critical value shows that SLOC is  $\leq 22$  for 70% of the Y-axis. The second (33) and the third (61) critical values represent 80% and 90% of the Y-axis respectively. In step 6, we can now find the range to define a method's size category: small ( $\text{SLOC} \leq 22$ ), medium ( $22 < \text{SLOC} \leq 33$ ), large ( $33 < \text{SLOC} \leq 61$ ), and very large ( $\text{SLOC} > 61$ ). Clearly, the results



(a) # revisions



(b) New additions



(c) Edit distance

Fig. 6: Cumulative distribution functions of percent of improvement in goodness of fit scores. All six metrics improve the prediction accuracy of the regression models when they are added with size. The two boxed areas in Figure (a) show why selecting only a few projects can be inaccurate for providing a generalizable observation about code metric usefulness. By selecting one boxed group only, we can underestimate (or overestimate) the effectiveness of code metrics.

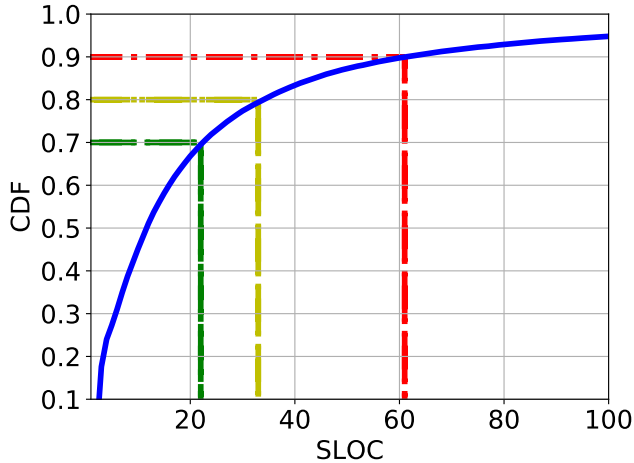
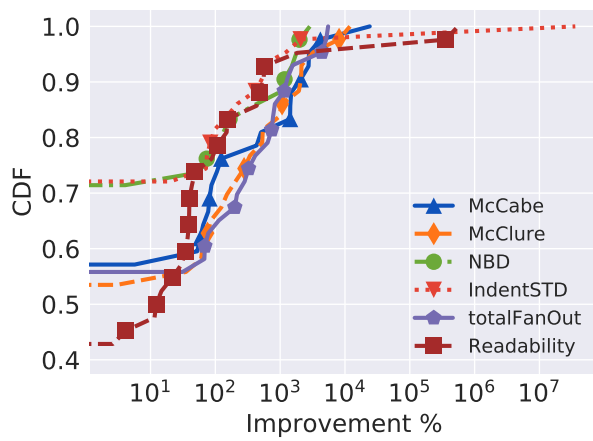


Fig. 7: Finding the critical values to determine small, medium, large, and very large methods, as proposed by Alves *et al.* [88].

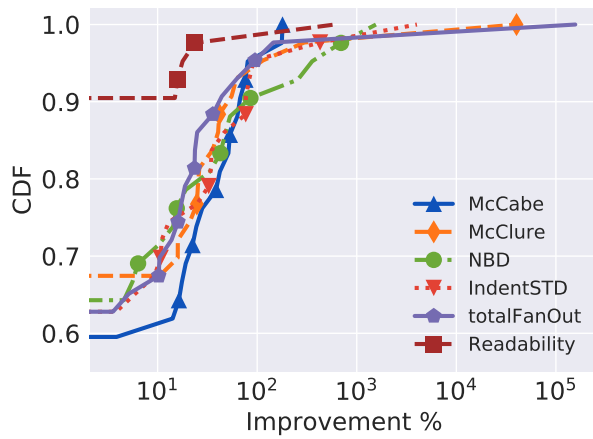
in Figure 6 are dominated by methods with  $\text{SLOC} \leq 33$  (80% of the Y-axis), limiting our understanding of metric performance for large methods. Here, we evaluate metrics for large and very large methods only ( $\text{SLOC} > 33$ )—we refer to both groups as “large” for simplicity.

Figure 8 shows that code metrics can be used to understand maintenance effort for large methods. However, their ranks in performance are not the same when compared with methods from all sizes (Figure 6, dominated by small SLOCs). We make the following observations while comparing Figure 6 and Figure 8: 1) Code metrics are useful for estimating maintenance effort for both large and small methods. 2) FanOut becomes one of the best maintenance predictors for large methods, although it was not among the best when all the methods were considered. This indicates that developers should be more careful about coupling (i.e., dependency on other methods) for large methods. 3) Readability effectiveness drops significantly (e.g., Figure 8(b), and 8(c)) when considering for large methods only. For edit distance for example (Figure 8(c)), Readability was able to improve the maintenance effort prediction accuracy only for 3 projects. 4) For edit distance, NBD and IndentSTD outperform all other metrics when the evaluation was dominated by small methods (Figure 6(c)). Surprisingly, their performance drops significantly for large methods. For example, NBD performs poorly for large methods with edit distance as the indicator, although it is the best for the same maintenance indicator when all methods are used. Although one may initially assume these observations as random noise, next we show that most of these performance variations are indeed explicable. We, therefore, need to consider the evaluation contexts to truly understand the usefulness of code metrics in estimating software maintenance effort.

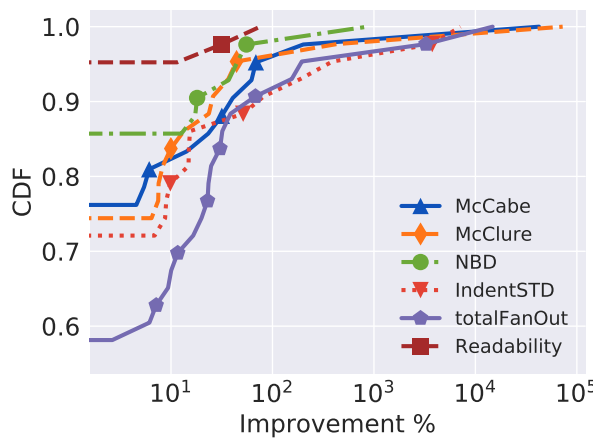
**Insight into the inconsistency of code metric performance:** Source code metrics will never model maintenance effort with 100% accuracy, because



(a) # revisions



(b) New additions



(c) Edit distance

Fig. 8: Percent of improvement in goodness of fit scores for large and very large methods only. Notice that all the metrics fail to improve the fitness score for a fraction of the projects. This is because many of these projects do not have enough large methods to produce regression coefficients with statistical significance ( $p$ -values are  $\geq 0.05$ ).

there are other factors that influence how a code component evolves over time: developer habits [89], application domain and platforms [50, 90], code clones [74], software architecture [91], and test code quality [18, 92]. Despite this difficulty, we take a deeper look into the following questions: i) Why is Readability’s performance so poor for large methods although it is excellent for all methods (i.e., dominated by small methods)? ii) Why does NBD perform the best for estimating edit distance when considered for all methods? And can we deduce a common phenomenon that explains the inconsistent performance of source code metrics?

**i) Readability:** Figure 6(a) shows that Readability is the best metric for estimating the number of revisions when the evaluation is dominated by a large number of small methods. According to the Wilcoxon rank-sum test, the performance distribution of Readability is statistically different from others. According to the Cliff’s Delta test, Readability’s performance has negligible effect size with IndentSTD, small effect size with NBD, and large effect size with all the others. Although its performance drops and becomes similar to the others when other maintenance indicators are used, Readability performs extremely poorly when evaluated for large methods only. For example, for edit distance and large methods, Readability is outperformed by all with large to small effect sizes. The scatter diagram in Figure 9 indicates that the Readability metric by Buse *et al.* [37] can only distinguish between how readable two methods are if the methods are small in size; to this model, all the large methods are similarly less readable. Interestingly, we find that the readability model was based on small code snippets only (maximum SLOC was 11 [93]). Evidently, the model does not scale for large methods. The graph shows that the variability in the readability measurements starts dwindling after SLOC 11 (first box), and almost diminishes after SLOC 33 (second bar). This clearly explains why the Readability metric performs poorly for the large methods.

**ii) Nested Block Depth:** Figure 6(c) suggests that NBD (and IndentSTD), outperforms all other metrics, when edit distance is the maintenance indicator (with statistically different performance distribution, and non-negligible effect sizes). Levenshtein edit distance [76] counts even the number of white spaces when it captures how many characters have been added, deleted, or edited to convert one method version into another. Although initially it seems like unimportant information to capture, it can actually indicate if a modification was done inside a nested block (by capturing the leading spaces) or outside. This is a strength of edit distance as a maintenance indicator, because modification inside nested blocks are considered more bug prone [53, 63]. This observation explains the superior performance of NBD and IndentSTD for this maintenance indicator; among the six code metrics, only these two capture nesting information of a method.

But why does NBD become the second worst metric for estimating the same edit distance when evaluated for large methods only? We observe that NBD lacks *variability* in large methods. Let us consider two large methods. One method has 2 nested loops, each with depth 3. The other one has one

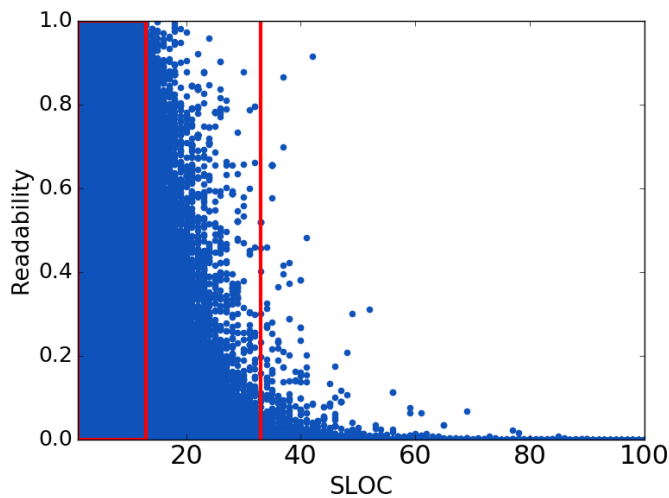


Fig. 9: Readability score against SLOC for all the methods in our dataset, based on the approach by Buse *et al.* [37]. SLOCs > 100 are discarded for graph readability.

nested loop only, but with the same depth. NBD as a metric fails to distinguish between these two methods— NBD is 3 for both. Some metrics such as FanOut, McCabe, and McClure do not have this limitation. Their value can increase monotonically with the increase in SLOC. These three metrics, in contrast to the others, never perform worse for large methods. In fact, they improved significantly for large methods. FanOut is the best predictor for edit distance (Figure 8(c)), whereas McCabe and McClure are the best predictor for new addition (Figure 8(b)).

To verify if our observation about measurement variability generalizes to all the metrics, let us consider Figure 10. Each metric measurement is represented by two box plots. The first shows the measurement distributions for all methods, and the second only for large methods. Except for the outliers, the second boxplot is always taller than the first one, because the measurements are naturally higher in large methods. When we compare the first boxplots for all metrics, they are not extremely different. However, the second boxplots are noticeably different from each other, and show that some metrics (FanOut, McCabe, and McClure) have much higher variability than others (e.g., NBD, and IndentSTD) when only the large methods are considered. Interestingly, a 1988 study by Weyuker [36] emphasized variability in measurements to be a desired property of a code metric, stating “a measure which rates all programs as equally complex is not really a measure”.

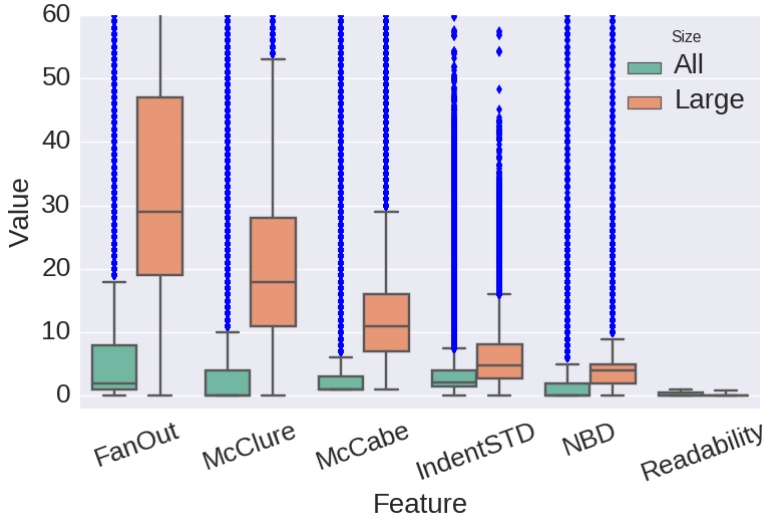


Fig. 10: Comparing variability in measurements between all methods and large methods for all the code metrics.

**Summary:** Code metrics are useful maintenance predictors, even after the size influence is completely neutralized. Their usefulness, however, depends greatly on the context in which they are applied. Some metrics are only good for particular maintenance indicators, and some only for small methods, because they reach a threshold and cannot discriminate further after the threshold measurement is reached.

## 6 Discussion

In this paper, we studied and reproduced the early contradictory claims about the relationship between maintenance impact and code metrics. We first focused on the methodological aspects that have influenced the outcomes of many previous studies. This investigation outlines some fundamental challenges that must be understood for accurately understanding a code metric’s validity. For example, if we can not neutralize the project’s difference in change evolution (Figure 1(b)), we should not rely on aggregated analysis. At the same time, enough projects should be considered to characterize outlying observations. While answering **RQ1**, and **RQ2** (*the confounding effect of size on the relationships between code metrics and maintenance effort*) we have established that the commonly used size normalization approach fails to neutralize the influence of size and should not be used in practice.

Using regression analysis for the size neutralized metrics evaluation (**RQ3**), we showed that we can use code metrics to prioritize our effort for reducing maintenance effort. This is encouraging for the research and development communities, because it refutes the claim that *size is the only valid code metric* [10,22,24], and it suggests that as a community we can continue researching new source code metrics that can provide greater insight into our software systems. The utility of a code metric, however, greatly depends on the evaluation



context in which it is applied. A metric, due to the lack of variability in the measurements, may become less useful when applied to large methods (**RQ4**). Additionally, a metric's performance is impacted by the underlying maintenance indicators used.

The varying performance of code metrics suggests that building context-aware maintenance models would be more effective than trying to derive a single generic model applicable to all systems. Software maintenance models have been studied for the last forty years [25], but considering their accuracy, there remains room for improvement [46]. Along with other existing recommendations, such as parameter optimization [94], we provide convincing evidence that the community should also focus on building ensemble [95] maintenance models instead of generic models that are commonly built (e.g., [46, 96, 97]). The envisioned approach is not to blindly apply a mixture of different machine learning algorithms (a form of ensemble modeling [98]), but to focus more on a mixture of models where each model is trained on a selected set of code metrics, a bounded method size, and a particular maintenance indicator. For example, while NBD and IndentSTD should be used for small methods and edit distance-based models, they should be excluded for other models that are better able to forecast maintenance indicators for large methods. This way we can build multiple base models based on the contextual code metrics evaluation, and can then aggregate the prediction of each base model that produces one final maintenance prediction for a given method [95].

### 6.1 Threats to Validity

Several threats may impact the observations in this study.

**Construct validity** is hampered by the maintenance indicators we used. Change proneness is not the only indicator of maintenance effort. Also, the indicators we used may not be sufficient to understand the true change proneness. Number of revisions, for example, can be impacted by the commit habits of contributors [5]. Some may commit more frequently than others. A less revised code, which is difficult to understand, and structurally complex to make a change, may require more effort than a more revised code. We mitigated this threat, at least to some extent, by using all the four change proneness indicators that we commonly found from the literature. The accuracy of collecting the complete change history in our measurements can be criticized, because we solely relied on CodeShovel [56]. Considering the run-time performance and accuracy, however, CodeShovel is the state-of-the-art tool for tracing method history.

**External validity** can be limited by the representativeness of the Java open source projects we used. Our results might not generalize to closed source projects. Also, we only focused on six selected code metrics. Although these six metrics were sufficient to demonstrate the problems and a potential solution for reliable evaluation of code metric performance, more code metrics should be investigated in the future.

**Internal validity** can be criticized by the statistical tests we used; however, these tests are well-established and well-recognized that seem appropriate for our context.

**Conclusion validity** can be hampered because of our dependence on regression analysis. We mitigated this threat by relying on regression coefficients that were statistically significant (i.e., p-value  $\leq 0.05$ ).

## 7 Conclusion & Future work

In this paper we set out to investigate whether code metrics can help us gain insight into maintenance effort, considering four code churn measures as maintenance effort indicators. While this question has been investigated before—and contradictory results have been presented—our approach was to reduce the level of granularity of our analysis to that of methods, and investigate the influence of size.

The key take away of our study is that *context is king*. Code metrics are useful predictors of maintenance effort, even after normalizing for size. However, their utility for predicting maintenance effort depends greatly on the context of how they are applied based on the type of maintenance indicators that are used and the size of the methods being examined.

This study presents a call-to-arms to the research community to investigate maintenance models that are context-aware, both in terms of method sizes and maintenance indicators. Implicitly, this is also a stringent warning for software engineering practitioners to not blindly follow metrics without taking context into account.

As our focus was on the previously unexplored method level granularity, we could not investigate the usefulness of some widely used class level code metrics (e.g., depth of inheritance). It would be interesting future work to see if the famous class level code metrics indeed help estimating future maintenance effort. We also plan to study the relationship between code metrics and bug-proneness with dedicated dataset reporting manually curated bugs.

## References

1. J. Börstler and B. Paech, “The role of method chains and comments in software readability and comprehension—an experiment,” *IEEE Transactions on Software Engineering*, vol. 42, no. 9, pp. 886–898, 2016.
2. D. Kafura and G. R. Reddy, “The use of software complexity metrics in software maintenance,” *IEEE Transactions on Software Engineering*, vol. SE-13, no. 3, pp. 335–343, 1987.
3. Y. Shin, A. Meneely, L. Williams, and J. A. Osborne, “Evaluating complexity, code churn, and developer activity metrics as indicators of software vulnerabilities,” *IEEE Transactions on Software Engineering*, vol. 37, no. 6, pp. 772–787, 2011.
4. Y. Zhou, B. Xu, and H. Leung, “On the ability of complexity metrics to predict fault-prone classes in object-oriented systems,” *Journal of Systems and Software*, vol. 83, no. 4, pp. 660 – 674, 2010.
5. D. Ståhl, A. Martini, and T. Mårtensson, “Big bangs and small pops: On critical cyclomatic complexity and developer integration behavior,” in *2019 IEEE/ACM 41st International Conference on Software Engineering: (ICSE-SEIP)*, 2019, pp. 81–90.
6. L. Cruz, R. Abreu, J. Grundy, L. Li, and X. Xia, “Do energy-oriented changes hinder maintainability?” in *2019 IEEE International Conference on Software Maintenance*

- and Evolution*, 2019, pp. 29–40.
7. M. Kondo, D. M. German, O. Mizuno, and E.-H. Choi, “The impact of context metrics on just-in-time defect prediction,” *Empirical software engineering*, vol. 25, no. 1, pp. 890–939, 2020.
  8. C. L. McClure, “A model for program complexity analysis,” in *Proceedings of the 3rd International Conference on Software Engineering*, 1978, p. 149–157.
  9. A. Tosun, A. Bener, B. Turhan, and T. Menzies, “Practical considerations in deploying statistical methods for defect prediction: A case study within the turkish telecommunications industry,” *Information and Software Technology*, vol. 52, no. 11, pp. 1242 – 1257, 2010.
  10. Y. Gil and G. Lalouche, “On the correlation between size and metric validity,” *Empirical Software Engineering*, vol. 22, no. 5, p. 2585–2611, Oct. 2017.
  11. T. J. McCabe, “A complexity measure,” *IEEE Transactions on Software Engineering*, vol. SE-2, no. 4, pp. 308–320, 1976.
  12. S. R. Chidamber and C. F. Kemerer, “A metrics suite for object oriented design,” *IEEE Transactions on Software Engineering*, vol. 20, no. 6, pp. 476–493, 1994.
  13. A. Lake and C. R. Cook, “Use of factor analysis to develop oop software complexity metrics,” USA, Tech. Rep., 1994.
  14. M. Shepperd, “A critique of cyclomatic complexity as a software metric,” *Software Engineering Journal*, vol. 3, no. 2, pp. 30–36, 1988.
  15. Y. Gil and G. Lalouche, “When do software complexity metrics mean nothing? – when examined out of context,” *Journal of Object Technology*, vol. 15, no. 1, pp. 2:1–25, Feb. 2016.
  16. J. Johnson, S. Lubo, N. Yedla, J. Aponte, and B. Sharif, “An empirical study assessing source code readability in comprehension,” in *2019 IEEE International Conference on Software Maintenance and Evolution*, 2019, pp. 513–523.
  17. D. Landman, A. Serebrenik, and J. Vinju, “Empirical analysis of the relationship between cc and sloc in a large corpus of java methods,” in *IEEE International Conference on Software Maintenance and Evolution*, 2014, pp. 221–230.
  18. D. Spadini, F. Palomba, A. Zaidman, M. Bruntink, and A. Bacchelli, “On the relation of test smells to software code quality,” in *2018 IEEE International Conference on Software Maintenance and Evolution*, 2018, pp. 1–12.
  19. R. K. Bandi, V. K. Vaishnavi, and D. E. Turk, “Predicting maintenance performance using object-oriented design complexity metrics,” *IEEE Transactions on Software Engineering*, vol. 29, no. 1, pp. 77–87, 2003.
  20. V. Antinyan, M. Staron, W. Meding, P. Österström, E. Wikstrom, J. Wrangler, A. Henriksson, and J. Hansson, “Identifying risky areas of software code in agile/lean software development: An industrial experience report,” in *IEEE Conference on Software Maintenance, Reengineering, and Reverse Engineering*, 2014, pp. 154–163.
  21. S. Scalabrino, G. Bavota, C. Vendome, M. Linares-Vásquez, D. Poshyanyk, and R. Oliveto, “Automatically assessing code understandability: How far are we?” in *32nd IEEE/ACM International Conference on Automated Software Engineering*, 2017, pp. 417–427.
  22. K. El Emam, S. Benlarbi, N. Goel, and S. N. Rai, “The confounding effect of class size on the validity of object-oriented metrics,” *IEEE Transactions on Software Engineering*, vol. 27, no. 7, pp. 630–650, 2001.
  23. D. I. K. Sjøberg, A. Yamashita, B. C. D. Anda, A. Mockus, and T. Dybå, “Quantifying the effect of code smells on maintenance effort,” *IEEE Transactions on Software Engineering*, vol. 39, no. 8, pp. 1144–1156, 2013.
  24. I. Herraiz, J. M. Gonzalez-Barahona, and G. Robles, “Towards a theoretical model for software growth,” in *Fourth International Workshop on Mining Software Repositories*, 2007, pp. 21–21.
  25. V. Lenarduzzi, A. Sillitti, and D. Taibi, “Analyzing forty years of software maintenance models,” in *International Conference on Software Engineering Companion (ICSE-C)*, 2017, pp. 146–148.
  26. U. Tiwari and S. Kumar, “Cyclomatic complexity metric for component based software,” *SIGSOFT Softw. Eng. Notes*, vol. 39, no. 1, p. 1–6, Feb. 2014.

27. M. A. Subandri and R. Sarno, "Cyclomatic complexity for determining product complexity level in cocomo ii," *Procedia Computer Science*, vol. 124, pp. 478 – 486, 2017, 4th Information Systems International Conference 2017, ISICO 2017, 6-8 November 2017, Bali, Indonesia.
28. D. Romano and M. Pinzger, "Using source code metrics to predict change-prone java interfaces," in *2011 27th IEEE International Conference on Software Maintenance*, 2011, pp. 303–312.
29. S. D. Suh and I. Neamtiu, "Studying software evolution for taming software complexity," in *21st Australian Software Engineering Conference*, 2010, pp. 3–12.
30. B. Robert, J. P. Correia, K. Schill, and J. Visser, "Standardized code quality benchmarking for improving software maintainability," *Software Quality Journal*, vol. 20, pp. 287–307, 2012.
31. Y. T. Chen, R. Gopinath, A. Tadakamalla, M. D. Ernst, R. Holmes, G. Fraser, P. Ammann, and R. Just, "Revisiting the relationship between fault detection, test adequacy criteria, and test set size," in *2020 35th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 2020, pp. 237–249.
32. C. Ebert, J. Cain, G. Antoniol, S. Counsell, and P. Laplante, "Cyclomatic complexity," *IEEE Software*, vol. 33, no. 6, pp. 27–29, 2016.
33. J. Pantiuchina, M. Lanza, and G. Bavota, "Improving code: The (mis) perception of quality metrics," in *IEEE International Conference on Software Maintenance and Evolution*, 2018, pp. 80–91.
34. B. Curtis, S. B. Sheppard, P. Milliman, M. A. Borst, and T. Love, "Measuring the psychological complexity of software maintenance tasks with the halstead and mccabe metrics," *IEEE Transactions on Software Engineering*, vol. SE-5, no. 2, pp. 96–104, 1979.
35. M. Alfadel, A. Kobilica, and J. Hassine, "Evaluation of halstead and cyclomatic complexity metrics in measuring defect density," in *2017 9th IEEE-GCC Conference and Exhibition*, 2017, pp. 1–9.
36. E. J. Weyuker, "Evaluating software complexity measures," *IEEE Transactions on Software Engineering*, vol. 14, no. 9, pp. 1357–1365, 1988.
37. R. P. L. Buse and W. R. Weimer, "Learning a metric for code readability," *IEEE Trans. Softw. Eng.*, vol. 36, no. 4, p. 546–558, Jul. 2010.
38. N. J. Abid, B. Sharif, N. Dragan, H. Alrasheed, and J. I. Maletic, "Developer reading behavior while summarizing java methods: Size and context matters," in *Proceedings of the 41st International Conference on Software Engineering*, 2019, p. 384–395.
39. J. Hofmeister, J. Siegmund, and D. V. Holt, "Shorter identifier names take longer to comprehend," in *IEEE 24th International Conference on Software Analysis, Evolution and Reengineering*, 2017, pp. 217–227.
40. S. Scalabrino, M. Linares-Vásquez, D. Poshyvanyk, and R. Oliveto, "Improving code readability models with textual features," in *IEEE 24th International Conference on Program Comprehension*, 2016, pp. 1–10.
41. V. Antinyan, M. Staron, and A. Sandberg, "Evaluating code complexity triggers, use of complexity measures and the influence of code complexity on maintenance time," *Empirical Softw. Engg.*, vol. 22, no. 6, p. 3057–3087, Dec. 2017.
42. J. Bauer, J. Siegmund, N. Peitek, J. C. Hofmeister, and S. Apel, "Indentation: simply a matter of style or support for program comprehension?" in *IEEE/ACM 27th International Conference on Program Comprehension*, 2019, pp. 154–164.
43. D. P. Darcy, C. F. Kemerer, S. A. Slaughter, and J. E. Tomayko, "The structural complexity of software an experimental test," *IEEE Transactions on Software Engineering*, vol. 31, no. 11, pp. 982–995, 2005.
44. J. M. Brittain, "Pitfalls of user research, and some neglected areas," *Social Science Information Studies*, vol. 2, no. 3, pp. 139–148, 1982.
45. L. Yu and A. Mishra, "An empirical study of lehman's law on software quality evolution," 2013.
46. L. Pascarella, F. Palomba, and A. Bacchelli, "On the performance of method-level bug prediction: A negative result," *Journal of Systems and Software*, vol. 161, 2020.
47. Q. Wang, X. Xia, D. Lo, and S. Li, "Why is my code change abandoned?" *Information and Software Technology*, vol. 110, pp. 108 – 120, 2019.

48. K. Herzig and A. Zeller, "The impact of tangled code changes," in *2013 10th Working Conference on Mining Software Repositories*, 2013, pp. 121–130.
49. D. Matter, A. Kuhn, and O. Nierstrasz, "Assigning bug reports using a vocabulary-based expertise model of developers," in *2009 6th IEEE International Working Conference on Mining Software Repositories*, 2009, pp. 131–140.
50. F. Zhang, A. Mockus, Y. Zou, F. Khomh, and A. E. Hassan, "How does context affect the distribution of software maintainability metrics?" in *IEEE International Conference on Software Maintenance*, 2013, pp. 350–359.
51. D. Radjenović, M. Heričko, R. Torkar, and A. Živković, "Software fault prediction metrics: A systematic literature review," *Information and Software Technology*, vol. 55, no. 8, pp. 1397 – 1418, 2013.
52. F. Palomba, A. Zaidman, R. Oliveto, and A. De Lucia, "An exploratory study on the relationship between changes and refactoring," in *Proceedings of the 25th International Conference on Program Comprehension*, 2017, p. 176–185.
53. A. Hindle, M. W. Godfrey, and R. C. Holt, "Reading beside the lines: Indentation as a proxy for complexity metric," in *16th IEEE International Conference on Program Comprehension*, 2008, pp. 133–142.
54. A. T. Ying, G. C. Murphy, R. Ng, and M. C. Chu-Carroll, "Predicting source code changes by mining change history," *IEEE transactions on Software Engineering*, vol. 30, no. 9, pp. 574–586, 2004.
55. Y. Higo, S. Hayashi, and S. Kusumoto, "On tracking java methods with git mechanisms," *Journal of Systems and Software*, vol. 165, p. 110571, 2020.
56. F. Grund, S. Chowdhury, N. C. Bradley, B. Hall, and R. Holmes, "Codeshovel: Constructing method-level source code histories," in *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*, 2021, pp. 1510–1522.
57. T. Menzies, J. Greenwald, and A. Frank, "Data mining static code attributes to learn defect predictors," *IEEE Transactions on Software Engineering*, vol. 33, no. 1, pp. 2–13, 2007.
58. E. Shihab, A. E. Hassan, B. Adams, and Z. M. Jiang, "An industrial study on the risk of software changes," in *Proceedings of the ACM SIGSOFT 20th International Symposium on the Foundations of Software Engineering*, 2012.
59. E. Giger, M. D'Ambros, M. Pinzger, and H. C. Gall, "Method-level bug prediction," in *Proceedings of the 2012 ACM-IEEE International Symposium on Empirical Software Engineering and Measurement*, 2012, pp. 171–180.
60. P. Ralph and E. Tempero, "Construct validity in software engineering research and software metrics," in *Proceedings of the 22nd International Conference on Evaluation and Assessment in Software Engineering 2018*, 2018, p. 13–23.
61. G. Sridhara, E. Hill, D. Muppaneni, L. Pollock, and K. Vijay-Shanker, "Towards automatically generating summary comments for java methods," in *Proceedings of the IEEE/ACM International Conference on Automated Software Engineering*, 2010, p. 43–52.
62. T. Zimmermann, R. Premraj, and A. Zeller, "Predicting defects for eclipse," in *Proceedings of the Third International Workshop on Predictor Models in Software Engineering*, 2007, p. 9.
63. M. R. Islam and M. F. Zibran, "How bugs are fixed: Exposing bug-fix patterns with edits and nesting levels," in *Proceedings of the 35th Annual ACM Symposium on Applied Computing*, 2020, pp. 1523–1531.
64. V. Antinyan, M. Staron, J. Derehag, M. Runsten, E. Wikström, W. Meding, A. Henriks-son, and J. Hansson, "Identifying complex functions: By investigating various aspects of code complexity," in *2015 Science and Information Conference (SAI)*, 2015, pp. 879–888.
65. R. Mo, Y. Cai, R. Kazman, L. Xiao, and Q. Feng, "Decoupling level: A new metric for architectural maintenance complexity," in *2016 IEEE/ACM 38th International Conference on Software Engineering*, 2016, pp. 499–510.
66. M. S. Rahman and C. K. Roy, "On the relationships between stability and bug-proneness of code clones: An empirical study," in *2017 IEEE 17th International Working Conference on Source Code Analysis and Manipulation (SCAM)*, 2017, pp. 131–140.

67. G. Catolino, F. Palomba, A. De Lucia, F. Ferrucci, and A. Zaidman, "Enhancing change prediction models using developer-related factors," *Journal of Systems and Software*, vol. 143, pp. 14–28, 2018.
68. A. Mocku and L. G. Votta, "Identifying reasons for software changes using historic databases," in *Proceedings 2000 International Conference on Software Maintenance*, 2000, pp. 120–130.
69. F. Khomh, M. Di Penta, Y.-G. Guéhéneuc, and G. Antoniol, "An exploratory study of the impact of antipatterns on class change- and fault-proneness," *Empirical software engineering : an international journal*, vol. 17, no. 3, pp. 243–275, 2012.
70. B. Ray, V. Hellendoorn, S. Godhane, Z. Tu, A. Bacchelli, and P. Devanbu, "On the "naturalness" of buggy code," in *Proceedings of the 38th International Conference on Software Engineering*, ser. ICSE '16, 2016, p. 428–439.
71. R. Moser, W. Pedrycz, and G. Succi, "Analysis of the reliability of a subset of change metrics for defect prediction," in *Proceedings of the Second ACM-IEEE International Symposium on Empirical Software Engineering and Measurement*, ser. ESEM '08, 2008, p. 309–311.
72. R. M. Bell, T. J. Ostrand, and E. J. Weyuker, "Does measuring code change improve fault prediction?" in *Proceedings of the 7th International Conference on Predictive Models in Software Engineering*, ser. Promise '11, 2011.
73. G. Bavota, M. Linares-Vásquez, C. E. Bernal-Cárdenas, M. Di Penta, R. Oliveto, and D. Poshyvanyk, "The impact of api change- and fault-proneness on the user ratings of android apps," *IEEE Transactions on Software Engineering*, vol. 41, no. 4, pp. 384–407, 2015.
74. A. Monden, D. Nakae, T. Kamiya, S. Sato, and K. Matsumoto, "Software quality analysis by code clones in industrial legacy software," in *Proceedings IEEE Symposium on Software Metrics*, 2002, pp. 87–94.
75. I. Scholtes, P. Mavrodiev, and F. Schweitzer, "From aristotle to ringelmann: a large-scale analysis of team productivity and coordination in open source software projects," *Empirical software engineering : an international journal*, vol. 21, no. 2, pp. 642–683, 2016.
76. V. I. Levenshtein, "Binary codes capable of correcting deletions, insertions, and reversals," in *Soviet physics doklady*, vol. 10, no. 8, 1966, pp. 707–710.
77. H. Hata, O. Mizuno, and T. Kikuno, "Historage: Fine-grained version control system for java," in *Proc. International Workshop on Principles of Software Evolution and ERCIM Workshop on Software Evolution*, 2011, pp. 96–100.
78. F. Grund, S. Chowdhury, N. C. Bradley, B. Hall, and R. Holmes, "Codeshovel: A reusable and available tool for extracting source code histories," in *2021 IEEE/ACM 43rd International Conference on Software Engineering: Companion Proceedings (ICSE-Companion)*, 2021, pp. 221–222.
79. M. M. Lehman, J. F. Ramil, P. D. Wernick, D. E. Perry, and W. M. Turski, "Metrics and laws of software evolution-the nineties view," in *International Software Metrics Symposium*, 1997, pp. 20–32.
80. H. C. Thode, *Testing for normality*. CRC press, 2002, vol. 164.
81. L. Inozemtseva and R. Holmes, "Coverage is not strongly correlated with test suite effectiveness," in *Proceedings of the 36th International Conference on Software Engineering*, 2014, p. 435–445.
82. S. Chowdhury, S. Borle, S. Romansky, and A. Hindle, "Greenscaler: training software energy models with automatic test generation," *Empirical software engineering : an international journal*, vol. 24, no. 4, pp. 1649–1692, 2019.
83. D. J. Sheskin, *Handbook of parametric and nonparametric statistical procedures*. CRC Press, 2020.
84. J. Romano, J. D. Kromrey, J. Coraggio, and J. Skowronek, "Appropriate statistics for ordinal level data: Should we really be using t-test and cohen'sd for evaluating group differences on the nsse and other surveys," in *annual meeting of the Florida Association of Institutional Research*, 2006, pp. 1–33.
85. R. Gopinath, C. Jensen, and A. Groce, "Code coverage for suite evaluation by developers," in *Proceedings of the 36th International Conference on Software Engineering*, 2014, pp. 72–82.

86. R. Just, D. Jalali, L. Inozemtseva, M. D. Ernst, R. Holmes, and G. Fraser, "Are mutants a valid substitute for real faults in software testing?" in *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering*, 2014, p. 654–665.
87. M. Papadakis, D. Shin, S. Yoo, and D.-H. Bae, "Are mutation scores correlated with real fault detection? a large scale empirical study on the relationship between mutants and real faults," in *Proceedings of the 40th International Conference on Software Engineering*, 2018, p. 537–548.
88. T. L. Alves, C. Ypma, and J. Visser, "Deriving metric thresholds from benchmark data," in *IEEE International Conference on Software Maintenance*, 2010, pp. 1–10.
89. A. Terceiro, L. R. Rios, and C. Chavez, "An empirical study on the structural complexity introduced by core and peripheral developers in free software projects," in *Brazilian Symposium on Software Engineering*, 2010, pp. 21–29.
90. M. Viggiano, J. Oliveira, E. Figueiredo, P. Jamshidi, and C. Kästner, "How do code changes evolve in different platforms? a mining-based investigation," in *2019 IEEE International Conference on Software Maintenance and Evolution*, 2019, pp. 218–222.
91. M. F. Aniche, C. Treude, A. Zaidman, A. van Deursen, and M. A. Gerosa, "SATT: tailoring code metric thresholds for different software architectures," in *16th IEEE International Working Conference on Source Code Analysis and Manipulation, 2016, Raleigh, NC, USA, October 2-3, 2016*, 2016, pp. 41–50.
92. D. Athanasiou, A. Nugroho, J. Visser, and A. Zaidman, "Test code quality and its relation to issue handling performance," *IEEE Trans. Software Eng.*, vol. 40, no. 11, pp. 1100–1125, 2014.
93. D. Posnett, A. Hindle, and P. Devanbu, "A simpler model of software readability," in *Proceedings of the 8th Working Conference on Mining Software Repositories*, 2011, pp. 73–82.
94. C. Tantithamthavorn, S. McIntosh, A. E. Hassan, and K. Matsumoto, "Automated parameter optimization of classification techniques for defect prediction models," in *IEEE/ACM 38th International Conference on Software Engineering*, 2016, pp. 321–332.
95. V. Kotu and B. Deshpande, "Chapter 2 - data mining process," in *Predictive Analytics and Data Mining*, V. Kotu and B. Deshpande, Eds. Boston: Morgan Kaufmann, 2015, pp. 17 – 36.
96. D. Coleman, D. Ash, B. Lowther, and P. Oman, "Using metrics to evaluate software system maintainability," *Computer*, vol. 27, no. 8, pp. 44–49, 1994.
97. N. Nagappan and T. Ball, "Use of relative code churn measures to predict system defect density," in *Proceedings. 27th International Conference on Software Engineering*, 2005, pp. 284–292.
98. H. Alsolai, M. Roper, and D. Nassar, "Predicting software maintainability in object-oriented systems using ensemble techniques," in *2018 IEEE International Conference on Software Maintenance and Evolution*, 2018, pp. 716–721.