

The 1 Quality Metrics of Test Suites in Test-Driven Designed Applications

[Alkaoud, Hessah Hassan]
Department of Computer Science
University of Colorado at Colorado Springs, USA
halkaoud@gmail.com

Kristen Walcott-Justice
Department of Computer Science
University of Colorado at Colorado Springs, USA
kjjustice@uccs.edu

ABSTRACT

As software engineering methodologies continue to progress, new techniques for writing and developing software have developed. One of these methods is known as Test-Driven Development (TDD). In TDD, tests are written first, prior to code development. Due to the nature of TDD, no code should be written without first having a test to execute it. Thus, in terms of code coverage, the quality of test suites written using TDD should be high. However, it is unclear if test suites developed using TDD exhibit high quality in terms of other metrics such as mutation score. Moreover, the association between coverage and mutation score is little studied.

reformat → too strong? Is it just that it is likely to be high??
In this work, we analyze applications written using TDD and more traditional techniques. Specifically, we demonstrate the quality of the associated test suites based on two quality metrics: 1) structure-based criterion represented by branch coverage and statement coverage, 2) fault-based criterion represented by mutation scores. We hypothesize that test suites with high test coverage will also have high mutation scores, and we especially intend to show this in the case of TDD applications. We found that Test-Driven Development is an efficient approach that improves the quality of the test code to cover more of the source code and also to reveal more faults based on mutation analysis.

I prefer to avoid using something that "we" actually did.
Keywords "we" unless it is

Test Driven Development, Test experimentation, Measurement, quality metrics, Mutation Score

1. INTRODUCTION

In software engineering, writing tests is vital to the development and maintenance processes because it verifies the correctness of the implementation of the software through program execution and analysis. During development, testing is a way to verify the completeness of the requirements. Testing also helps to ensure that the software has achieved

Here you could probably write that "The results reveal that ..."

It is probably simpler to write that "This paper analyzes ...". Please rephrase so that the sentence does not use etc. and the ellipse.

an appropriate level of quality. In maintenance, software testing is important when adding new functionality or refactoring to check the correctness of the changes. Moreover, testing is important when verifying that the new feature does not cause the program to regress during maintenance work. Due to this importance, new methods of developing programs and tests together have been devised. One of these techniques is known as Test-Driven Development (TDD). TDD practices require that tests be written prior to the writing of code. In this methodology, no code should be written without first having a test case to cover such code. Thus, the quality of the code as measured using coverage should be high for TDD applications. However, there is little research in comparing code coverage of such applications to other code quality metrics.

I would write "the test cases"
Traditionally, tests are written in an ad hoc fashion after the code is completed by the developers, in tandem with development by a QA team, or, in TDD, tests are written prior to any code. The TDD process is also called test-first or red-green-refactoring. Today TDD is being widely adopted in industry, including large software firms such as Microsoft and IBM [27]. TDD also gained in popularity with the introduction of the eXtreme Programming XP methodology [7], and it is sometimes used as a stand-alone approach for software engineering tasks such as adding features to legacy code [13]. In TDD, the software is developed in short iterative cycles and in each cycle, testing is repeated for each testable component of the program. The pattern of the cycles can be explained as the following:

- Write the test and run it to see that it fails (Red).
- If the test fails, implement the code that will cause the test to run and pass the test (Green).
- Alter the design of both the test and the code (Refactoring).

This term has not been defined.

Refactoring is a process of altering the code for the purpose of enhancing its design, readability, maintainability, etc... without violating existing code functionality [14]

Also, maybe comment on the size of the programs & test suites.

Once a test suite is written and complete, it is also important to evaluate the quality of the test suite. Many program-based adequacy criterion have been developed to determine whether or not a test suite adequately tests the application [33]. Program-based adequacy criterion exploit

would it be useful if the abstract commented on the real-world nature of the programs?

→ This could be plural & thus I suggest using "criteria".

both structure-based criterion and fault-based criterion. A structure-based criterion requires the creation of a test suite that solely requires the exercising of certain control structures and variables within program [21]. Fault-based test adequacy criterion on the other hand, attempts to ensure that the program does not contain the types of faults that are commonly introduced into software systems by programmers [12, 34]. One of the most common types of test quality evaluation in both industry and academic work is based on structurally-based criterion which is commonly analyzed in terms of statement or branch coverage [33].

→ criteria... attempt?

When developing programs using TDD, code is only written when a test shows that it is needed. In theory, no production code in TDD should be developed unless it has at least one associated test. Therefore, the code coverage of TDD-based tests should be high for their associated applications ~~by nature~~. Both the literature and practice indicate that the use of TDD yields several benefits. Among others, TDD leads to improved test coverage [5]. However, statement coverage and branch coverage have been shown to not necessarily give firm evidence of high quality test suites.

→ by definition, or, by nature,

In recent years, mutation testing/scores have been used as another method for evaluating test suite quality. Mutation testing is a fault-based technique which measures the fault-finding effectiveness of test suites on the basis of induced faults [12, 17]. Mutation testing is a well-known technique to design a new software tests or to evaluate the quality of existing software tests and test suites. The idea of using mutants to measure test suite adequacy was originally proposed by DeMillo et al. [12]. Mutation testing involves seeding faults in the source program. Each altered version of the source is called a mutant. When a test reveals the mutant then the mutant said to be killed. The ratio of killed mutants/generated mutants is known as the mutation score. In mutation testing, in a simple statement such as `if (a < b)`, the `<` sign will be replaced with all the possible relational operators such as `>`, `<=`, `>=`, `==`, `!=`. The use of mutation operators that is used by tools yields results in the empirical assessment of quality for current testing techniques [3].

→ is the / needed.

Given the rise in popularity of TDD approaches, the question is whether TDD ~~paper~~ promotes high quality test suites or not? In this work, we learn if TDD methodology promotes higher test suite quality than traditional test design technology. To show this, we analyze and compare three quality metrics: branch coverage, statement coverage, and mutation score when applied to two sets of programs 1) ones that were designed using the TDD methodology and 2) programs designed using standard testing techniques. In this document, the term standard testing techniques refers to any tests that were developed using methodologies that are not TDD oriented. We first analyze the size and complexity of the programs. Then we demonstrate the statement and branch coverage quality of the existing test suites. Next, we apply a mutation testing tool to determine the mutation score of the test suite. Finally, we analyze the association between test coverage and mutation scores in estimating overall test suite quality.

"analyze"

Our results show that test-driven development has a higher quality of test than the other development approaches. We

→ I'm not sure that we actually need to put the 1) and 2) in the paper here.

→ It is better to use bullet points.

demonstrate a relationship between the mutation score and traditional coverage metrics, statement coverage and branch coverage, and that the correlation with branch coverage is stronger. A percentage of 54% of the programs reflects a relationship with statement coverage and more as 62% reflects a relation with branch coverage.

→ not needed → This is hard to parse - rephrase?

In summary, the main contributions of this research are as follows:

→ Tighten to fit on a single line?

- An examination of the quality of TDD test suites based on two quality metrics, code coverage and mutation score.
- An empirical evaluation of the quality of TDD test suites considering mutation and coverage scores.
- A discussion of the observed trends and comparison of the two quality metrics: coverage and mutation.

2. BACKGROUND

In this section, we discuss background regarding test-driven development practices and two program-based criteria that are used for evaluating test suite quality: code coverage and mutation scores.

2.1 TDD Practices

Test-driven development is a software development process that is based on repeating a very short development cycle. There are several variations of test-driven development, although the main method is called Test-First or testing via a red-green-refactoring cycle. Refactoring is a process of altering the code for the purpose of enhancing its design, readability, maintainability, etc. without violating existing code functionality [14]. The TDD process is performed in the following six steps:

- can you avoid the repeated use of "tech"?
- Add test code for a new functionality or a specific desired improvement. → I cannot easily understand the meaning of this sentence.
 - Run the test code and see the test fail (red).
 - Implement the production code according to the test.
 - Run all tests again and make sure they pass (green).
 - Repeat 3, 4 until the test passes.
 - Refactor the production code and the test. → There are no numbers.
 - Repeat steps 1 through 7 to program development completion.

→ Yes, but we might need to explain better how a test reveals the mutant.

Test-driven-development implies quality benefits to the software. It is true that TDD requires more test-based code than without TDD because of the thoroughness of the unit test code, but the large numbers of tests also contributes to limiting the number of defects in the code. The frequent and early nature of the TDD testing helps to reveal defects early in the development process, reducing the number of bugs that may be observed later in time. Moreover, TDD tests that are written prior to code development and are run repeatedly are more likely to detect any regression of the code's behavior. This allows developers to discover problems that may arise when a later change in the implementation unexpectedly affects other functionality.

→ best word ?? maybe "investigate"

In software engineering, behavior-driven development BDD is another software development process that is built based

→ We should consistently use the word "paper" in the final paper submission.

Points: (1) we may want to mention this as an example of future work.
(2) Do these results also apply to BDD??

on test-driven development (TDD). The original developer of BDD, Dan North [28], came up with the notion of BDD because he was dissatisfied with the lack of any specification within TDD of what should be tested and how. Astels notes, "Behavior-Driven Development is what you were doing already if you're doing Test-Driven Development very well [6]." In this work, we only focus on TDD-specific applications.

2.2 Value of Coverage *glad to read this content in the paper*

Code coverage metrics are common for evaluating the quality of test cases. Coverage metrics check how many elements of the source code have been *executed* by the test. Code coverage falls under structurally-based criterion that requires the creation of a test suite that solely requires the exercising of certain control structures, statements, or variables within the program[21]. There are many coverage metrics types that have been used and proposed, some of which include statement coverage, branch coverage, and condition coverage. Branch coverage criteria and statement coverage criteria are the two most popular examples of structural adequacy criteria that are used in existing coverage analysis tools and by industry[33].

Statement Coverage metrics require executing each statement in the program at least once during testing. It can be defined as follows:

Definition: the coverage of all nodes in the flow graph [34].

Branch Coverage metrics require the execution of each control transfer (true/false) of each branch (e.g. if statements, loops) in the program under test at least once. Branch coverage can be defined as follows:

Definition: the coverage of all edges in the flow graph [34].

To distinguish the differences between statement and branch coverage, consider the code below:

```
public int returnInput(int input, boolean condition1,  
boolean condition2, boolean condition3)  
int x = input;  
int y = 0;  
if (condition1)  
x++;  
if (condition2)  
x-;  
if (condition3)  
y=x;  
return y;
```

This should be formatted properly and placed inside of a figure.

with the following associated test case:

Test case 1: shouldReturnInput(x, true, true, true)

In this case, statement coverage would be 100% as every statement would be executed. Test case 1 will guarantee that each statement in the code will be exercised. However, the branch coverage is only 50% because branch coverage requires checking that both edges of each branch execute. To improve the branch coverage to 100%, a second test case could be added as follows:

Test case 2: shouldReturnInput(x, false, false, false)

Overall, I think that some of the content on this page can be removed or condensed. For instance, we do not need to define statement coverage formally.

From these examples of test cases, we observe that branch coverage implies statement coverage, because exercising every branch leads to executing every statement [11]. Branch coverage is a good indicator of test quality, but still branch coverage is not enough to truly estimate test quality because it does not imply the execution of all possible paths. The meaning will be explained in this example:

This simple method has two tests, and the branch coverage calculation reports 100% branch coverage even though there is missing a test when $i == 0$.

```
public static String foo(int i)  
if (i <= 0)  
return "foo";  
else  
return "bar";
```

@Test

```
public void shouldReturnBarWhenGiven1()  
assertEquals("bar", foo(1));
```

@Test

```
public void shouldReturnFooWhenGivenMinus1()  
assertEquals("foo", foo(-1));
```

Thus, high code coverage does not necessarily imply full coverage and testing overall.

There are papers that empirically support this claim

2.3 Mutation Score

Mutation analysis is an alternative approach for determining code quality. Mutation analysis was originally introduced in [12, 8]. Mutation testing is a fault-based technique which measures the fault-finding effectiveness of test suites on the basis of introduced faults[12, 17]. Fault-based test adequacy criterion attempt to ensure that the program does not contain the types of faults that are commonly introduced into software systems by programmers [12]. Mutation testing involves modifying a program's source code in small ways in order to seed in faults [29]. Each modified version of the program is called a *mutant*, and each mutant includes only one potential fault. Fault planting is based on a set of mutation operators. For example, mutation operators may include:

reformat

- Statement scratching.
- Replace the boolean sub-expression with true and false.
- Replace the arithmetic operation with another, e.g. $+$ with $*$, $-$ and $/$.
- Replace the boolean relation with another, e.g. $>$ with \geq , \equiv and \leq .
- Replace the variables with another declared variable in the same scope that has the same type.

part of the paper to score correctly

How do you handle these mutants??

In the test execution process, if the test revealed the fault then it is said that the test *killed* the mutant. *Equivalent* mutants are mutant programs that are functionally equivalent to the original program and therefore, cannot be killed by any test case.

I don't think that a comma is needed,

Relating the total number of the killed mutants to the total number of generated mutant is a powerful way to measure the quality of the test suites. The number of mutants killed given the number of mutants create results in the calculation

see [*] at bottom.

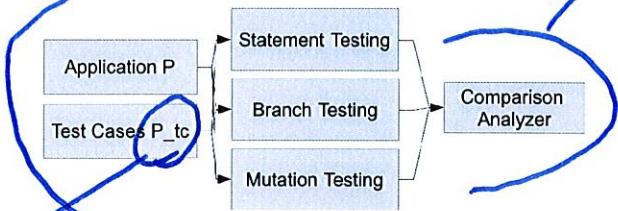


Figure 1: The testing framework.
correct phrase?

a mutation score (MS), which is a testing criterion to measure the effectiveness or ability of a test suite to detect faults.

$$\text{mutation_score} = \frac{\text{number_of_mutants_killed}}{\text{total_number_of_mutants}} \quad (1)$$

The mutation score determines whether the test is good enough to detect seeded faults. Mutation testing is not meant as a replacement for code coverage; rather, it is a complementary approach that is useful for finding code that is covered by tests based on coverage metrics but may still not be effectively tested [24]. While traditional code coverage measures how much of code is executed by the test suite, the mutation score measures the ability of tests running on developed code to reveal random faults. The simplest example of the usefulness of mutation adequacy is this conditional logic statement: `if (a < b)`. It is useful to include mutation operators that manipulate the relational operators found within a program. Thus, the '`<`' operator may be replaced with other possible relational operators such as `(>, <=, >=, ==, !=)` [16]. Ideally, *mutation* operators should produce the kinds of mutant programs that software engineers are most likely to create.

This does not need to be in italics.

Note that some mutants may not actually introduce faults, but rather create an equivalent program given the feasible paths in the program. In this work, we assume that mutants may or may not be equivalent to the original program, given the tools we use. Thus, a mutation score of 100% may not be possible for all programs. If equivalents could be easily assessed, a mutation score of 100% would be the goal. However, in this work, we aim for "high" mutation scores, where "high" will be judged relative to the mutation scores of other test suites and programs.

Again, "This paper compares...

3. ANALYZING TDD APPROACHES

In this research, we compare the statement coverage, branch coverage, and mutation scores of programs that have been developed using TDD and other practices. Our process can be seen in Figure 1. Programs and their existing test suites are executed using a coverage analyzer tool and a mutation analyzer tool. Given those results, we compare the coverage and mutation scores for different types of applications.

* Did Frank have work like this??
We expect that TDD approaches will exhibit high code coverage in terms of both statement and branch coverage, given the practices of TDD. However, there is little known about the correlation between coverage and mutation score. Applications that were developed using other methods are analyzed with the same tools. We expect that coverage and mutation scores will be lower than those developed using TDD, especially for applications that are not developed within a

What is the output of this process??

large community or are not widely used.

Reformat.

4. EMPIRICAL EVALUATION

In this research, we compare the quality of TDD test suites to other those developed using more standard testing techniques. First, we observe the trends of quality in test suites that are written using TDD methodology. Next, we observe the trends of quality when generating test suites using other development techniques. Then, we compare the quality of TDD test suites to other test suites developed using techniques and observe the overall impact of code coverage criteria and mutation score criteria. In summary, we analyze:

- how did you handle equivalent mutants??
- The relationship between coverage and mutation scores when applied to programs developed using TDD
 - The relationship between coverage and mutation scores when applied to programs developed using non-TDD methods *Try to reformat for a better fit.*
 - Potential relationships between coverage and mutation scores given program properties

Better to use notation than under-scores?

Overall, we examine how traditional software testing methodologies compare to newer testing methodologies, specifically test-driven-development (TDD) with regard to test suite quality. We answer three key questions *research questions*.

Q1: What is the impact of TDD on other test quality metrics? *Rewrite to fit on one line?*

In the way the TDD process is designed, it directly implies good coverage, as described in Sections 2. However, there is a little understanding in current research as to how quality can be measured other than coverage. Another potential quality metric includes mutation metrics and, more specifically, mutation score. If we know that the code coverage is high in TDD applications, does that imply that the mutation score is also high, nearing 100%?

We have also observed that TDD applications claim high "code coverage." In most work, it is unclear if code coverage refers to statement coverage, a lesser metric, or to branch coverage. Thus, we also analyze if TDD test suites are creating high quality test suites in terms of both statement and branch coverage.

Q2: What are the differences in the quality of test suites between TDD and other testing methodologies?

Once both quality metrics, coverage and mutation score, are taken into consideration, we next ask what qualities traditional test suites exhibit. Traditional test suites are often developed in ad hoc fashions. Some may develop the tests in tandem with the writing of code, or the code may be written and then passed off to an external (away from the developer) Quality Assurance (QA) team for evaluation. QA teams often test the code in an ad hoc or behaviorally or integration testing driven way. For example, "if I click this button, then I should arrive at this view and see these results." Tests written in this way account little for individual units of code or subtle integrations between components. We examine traditionally created test suites for applications and evaluate the code coverage and mutation scores for each. We hypothesize that the coverage and mutation scores will be lower than that of applications generated using TDD.

Q3: What are the trends in test quality metrics?

There are also some debatable questions that are related

[*] I really like these simple graphs. I wonder if you could redraw this using PGF/TikZ or some other tool that will help improve the diagram's quality.

Is there any way that you can know for sure that a given Java project was not developed w/ TDD? Look in the version control repository??

Better to put something leave than blank!

I think that this would be better expressed using a question & then checkmarks.

Benchmarks			Test					Program properties		
name	version	group	Statement coverage	Branch coverage	Mutation score	SLOC	Test-LOC / Source-LOC	CCN		
Tyburn	1.1	TDD	92	88	86	554	1.2	1.33		
JBehave-core	3.9	TDD	88	81	100	9726	0.8	1.57		
Helium	0.1	TDD	94	99	92	1018	1.7	1.95		
Trove	3.0.3	Non-TDD	7	7	59	2080	5.0	2.24		
Commons-lang	3.2	Non-TDD	94	90	75	15825	1.8	3.31		
Commons-io	2.4	Non-TDD	86	85	82	6037	2.0	2.52		
Jdom	2.0.0	Non-TDD	92	86	77	12038	1.4	2.58		
Numerics4j	1.2	Non-TDD	99	97	81	1836	1.5	2.27		
lavalamp	1.0	Non-TDD	98	35	72	1039	1.3	1.28		
netweaver		Non-TDD	16	2	15	17953	0.7	2.7		
jaxen	1.1.6	Non-TDD	79	59	50.41	8428	0.8	3.21		

Table 1: Test applications and statistics.

I would call these "subject programs" or "case study Applications".

In our experiment, we have incorporated three TDD projects. These benchmarks include Tyburn, JBehave and Helium. We also evaluate several non-TDD projects including Trove, Commons Lang, JDOM, Commons IO, numerics4j, and jaxen along with netweaver and lavalamp from SF100.

Is this how the word is normally written?

4.1.3 Metrics

We measure branch coverage, statement coverage, and mutation score on all projects. Each of these metrics is described in Section 2. Two main tools are used for these measurements: Jacoco and MAJOR. Jacoco [18] is an open source coverage library for Java, which has been created by the EclEmma team. Jacoco reports code coverage analysis (e.g. line, branch, instruction coverage) from the bytecode. MAJOR [20] is a mutation testing and score tool developed by René Just. The Major mutation framework enables fundamental research on mutation testing as well as efficient mutation analysis of large software systems. The mutation framework provides the following three components.

- Compiler-integrated mutator.
- Mutation analysis back-end for JUnit tests.
- Domain specific language to configure the mutation process in detail.

This word does not seem to fit here.

For the final paper, we do not need to give the name of the person who created the tool.

We also calculate program and test attribute measurements to help us recognize the distinct trends in the code projects. We measure project characteristics including: number of non-commented source code lines (SLOC), complexity of the source code, and test lines of code per source line of code (test-LOC/source-LOC). These calculations are performed using JavaNCSS. JavaNCSS [22] is a command-line tool that can be used to count the non-commented lines of source code, as well as the McCabe cyclomatic complexity of the source code. The cyclomatic complexity [25] metric for software systems is adapted from the classical graph theoretical cyclomatic number and can be defined as the number of linearly independent paths in a program. Larger numbers of lines of code and increased complexity can imply poor design and greater challenges to thorough testing.

So, does the SLOC include the test cases?

4.2 Results

We evaluate the source lines of code (SLOC), test to source lines of code (test-LOC/source-LOC), and the average complexity of production code for all projects. Also, we evaluate the statement coverage, branch code coverage, and mutation scores of all test suites of our selected benchmark programs,

① Can you clearly state that these values are percentages?

② There should be a footnote for the table to indicate the name of the tool used to calc.

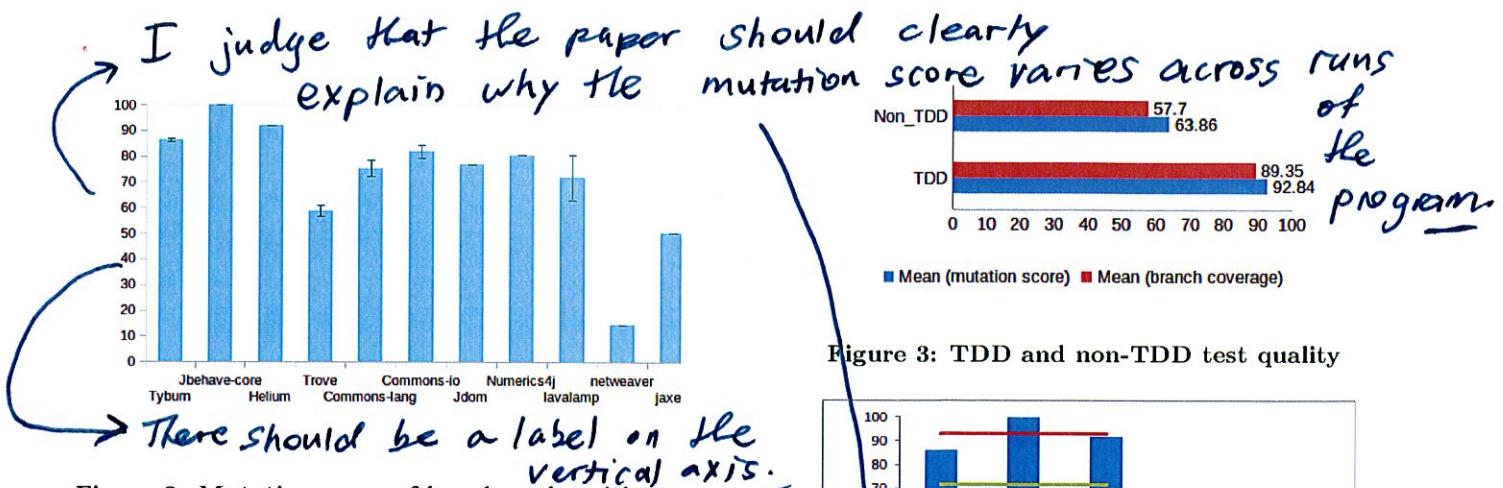


Figure 2: Mutation score of benchmarks with standard deviation

Group	n	Mean	Standard Deviation
TDD	3	89	9
Non-TDD	8	58	38

Table 2: Means of Branch coverage

Group	n	Mean	Standard Deviation
TDD	3	93	7
Non-TDD	8	64	23

Table 3: Means of Mutation score

as can be seen in Table 1.

Reformat.

4.2.1 TDD vs Traditional

In this section we aim to answer these two questions:

Q1: What is the impact of TDD on other test quality metrics?

Q2: What are the differences in the quality of test suites between TDD and other testing methodologies?

We measured the mutation score based on five executions for each program. Then we count the mode of the scores to represent the mutation scores. We also counted the standard deviation of the five executions. In Figure 2, the standard deviation is represented by the error bars for each program showing the range of the mutation scores.

Next we calculate the mean and the standard deviation of the two groups in terms of both quality metrics, branch coverage and mutation score. The results can be seen in Table 2.

In Figure 3, we compare the two groups' arithmetic means. The TDD group of programs show higher overall test quality than non-TDD programs when considering both mutation score and branch coverage. While the branch coverage mean of the non-TDD group is 57.7, the TDD group's mean is 89.35. Similarly, in terms of mutation score, the mean of the mutation score for TDD group is higher than the non-TDD group. The TDD programs achieve an average mutation score of 92.84 compared to 63.68 for non-TDD applications. Also, in Figure 4, we visualize a comparison between the arithmetic mean of the branch coverage and mutation score of the TDD group in comparison to non-TDD applications. The TDD applications show coverage and mutation scores

"visually compares"

Figure 3: TDD and non-TDD test quality

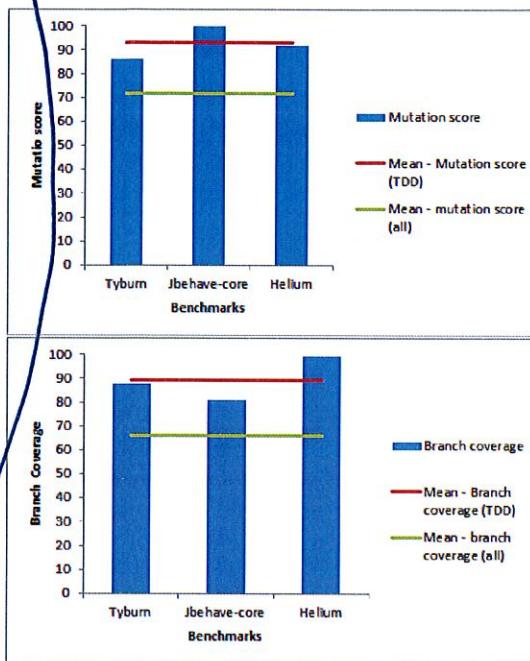


Figure 4: Mutation scores and Branch coverage scores for TDD applications.

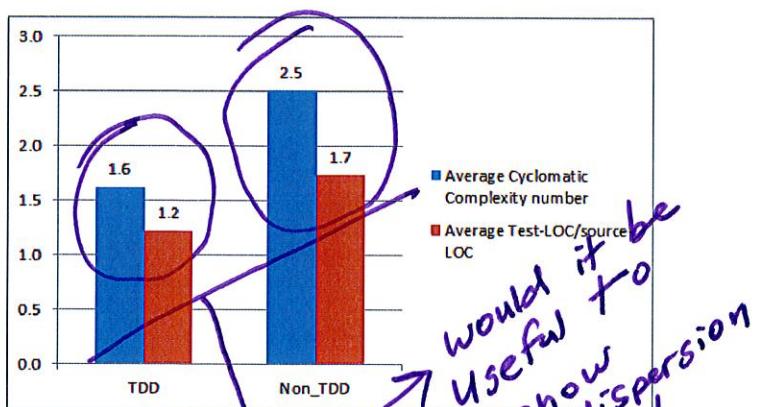


Figure 5: Comparing TDD and non-TDD applications given Test lines of code vs source lines of code and code cyclomatic complexity

above the average of all non-TDD benchmarks.

In Figure 5, we calculated the arithmetic mean of TDD and non-TDD applications with regard to test-line-of-code per

Reformat.

"well, the calc. is not in fig 5- it just visualizes the calc., right?"

I think that we should probably cite this book here!

In the book called "The Statistical Sleuth" the authors say that $R^2 = 7.5$ is good in the social sciences (which our work is somewhat like :)

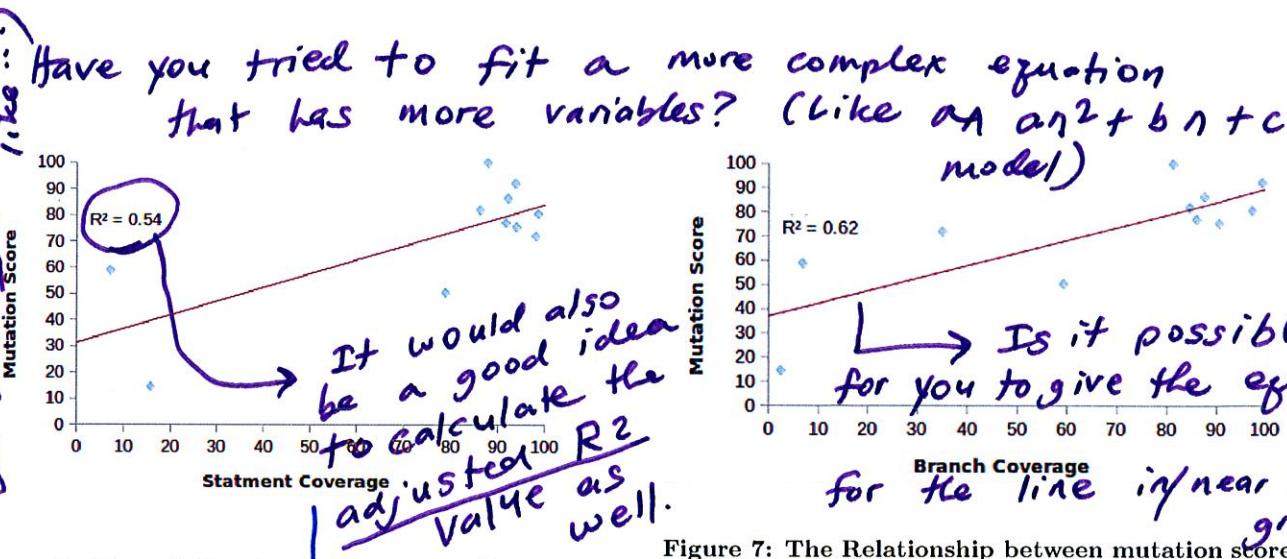


Figure 6: The relationship between mutation score and statement coverage.

source-line-of-code (Test-LOC / Source/LOC) and the average cyclomatic complexity for the functions of the source code (CCN). TDD programs tend to have lower CCN and tests. We believe this is because the TDD programs are better refactored to lower the complexity of the functions in the source code as well as refactoring the test code to remove potential duplication.

"anticipate that this is

4.2.2 Mutation Score versus Coverage due to

Our results show that there is a relationship between mutation score with both traditional coverage metrics: statement coverage and branch coverage.

As shown in Figure 6, there is a relation between the statement coverage and the mutation score. It seems the mutation score is dependent on the branch coverage, and it can be predicted based on the branch coverage, the independent variable. As visualized in the figure, the rounded R-squared value of 0.54 indicates an existent relationship between the branch coverage and the mutation score. In other words, almost 54% of the variation in branch coverage can be explained by the variation in mutation score.

While there is a correlation between mutation score and statement coverage, the correlation between mutation score and branch coverage is stronger. Figure 7 reflects a relation where the coefficient of determination is 0.62. In other words, the programs that had high branch coverage also achieved a high mutation score, as seen in Figure 7. The resulting formula is:

$$f(x) = 0.52x + 37.1 \quad (2)$$

Too much extra space.

Despite the fact that 62% of the programs can be explained by the relationship between mutation score and branch coverage, 36% of the programs in our study that cannot be explained with regard to this relation. The correlation between statement coverage cannot be explained by the observed branch and mutation score correlation. To illustrate, in Lavalamp branch coverage scores were 35% whereas the mutation score reaches 75%. We believe this was because the statement coverage of this program was as high as 98%. Some other programs cannot be explained by either

It might also be a good idea to do hypothesis testing with the Mann-Whitney U-Test and effect size calculation (i.e., A_{12}) — see an ICSE paper by Arcuri for details.

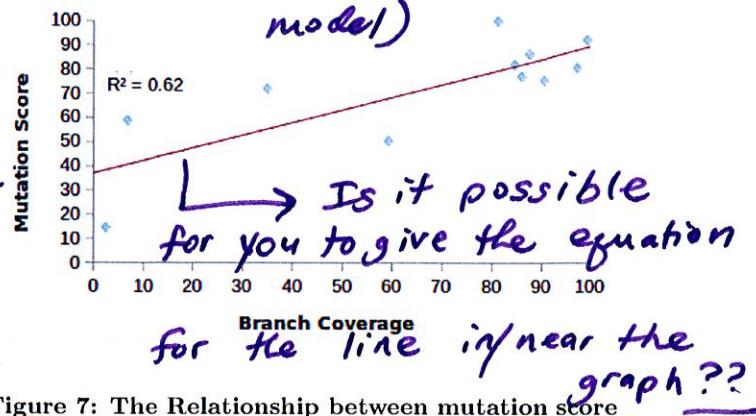


Figure 7: The Relationship between mutation score and branch coverage

Cody Kinner knows how to make

by branch or statement coverage. For example, in common-lang the branch coverage of common-lang is equal to 90% and the statement coverage is 94%. However, its mutation score is 75%.

these graphs and do this analysis in R.

5. DISCUSSION

In TDD, we found a trend between test-LOC per Source-Line of Code with high branch coverage and statement coverage. In other words, as more tests are written in TDD, branch and the statement coverage improve. The R-squared value is a linear relationship of 0.89 corresponding to statement coverage and 0.99 to branch coverage. This indicates that most of the tests in TDD are written to improve the branch coverage. It is known that executing every branch executes every statement as well unless there is a break or a jump [18]. Therefore, as a consequence of the high branch coverage in TDD programs, the statement coverage and the mutation score of TDD programs is also high.

MAJOR SUPPORTS THIS.

Mutation score metrics have a relationship not only with branch coverage but also with statement coverage. Both statement and branch coverage will complimentarily predict mutation score. Looking at Lavalamp as an example, the statement coverage is 98%, the branch coverage is 35% and the mutation score is 72%. Despite the low branch coverage of "Lavalamp" which was 35%, the mutation score of Lavalamp is 72%, which is not low. This could at first glance bring up a question about the relationship between the mutation score and the branch coverage. However, we show in our results that the mutation score can be predicted from both branch and statement coverage in most cases. In this case, the non-killed mutants of "Lavalamp" were due to the uncovered branches that haven't tested.

Our results summarize that writing more test cases does not necessarily ensure good test suite quality on both metrics. We found no relation between test-LOC/source-LOC. The coefficient of determination of Test-LOC with the mutation score range was 0.02 and was 0.05 with branch coverage. Writing more tests could give little extra value, such as in trove where the TLOC/SLOC is the highest among the projects as seen in Table 1, but the tests did not reflect a very high mutation score or coverage because most of trove's tests were written for one or two modules of the source, leaving

In some cases, it might be a good idea to look at

mutants that cannot be killed

extra space needed, "we're not"

better to use the same font convention

I think that using it is prob. the best.

Reformat so flat you do not have a single line at top or bottom of a page.

expensive, code coverage-based analyses.

8. CONCLUSION AND FUTURE WORK

In this paper, we performed a test suite quality based evaluation for the TDD testing process. We demonstrate the effect of TDD practices on test suite quality and compare that to other testing methodologies. Test suite quality is measured based on branch and statement coverage and compared to mutation scores. Finally, we discuss the correlation between code coverage quality and mutation scores based on our observations.

Reformat.

more "than" performed.
Also, report on

We have quantitatively evaluated Test-Driven-Development by applying test quality metrics to three open source projects. Our study gives evidence that the Test-driven development is indeed an effective approach for writing high quality test suites. For the TDD applications under consideration, Tyburn exhibited the lowest mutation score at 86% whereas Jbehave had a score of 100%. On the other hand, the mutation scores of non-TDD programs ranged from 15% in Netweaver to 82% for commons-io.

We also found that TDD programs tend to focus test writing based on code branches. We found a linear relationship between first test-line-of-code per source-line-of-code in TDD programs compared to branch coverage. The coefficient of determination R^2 of this relation is 0.99 where the maximum of value of R^2 is 1.

Finally, we discuss the correlation between code coverage quality and mutation scores based on our observations. We have discovered that the mutation score often can be predicted from both traditional metrics, statement and branch coverage. Both metrics, when examined together, indicate a higher likelihood to kill mutants. A percentage of 64% of the programs mutation scores can be anticipated based on the branch coverage number of that program. A lower percentage of 54% of the programs' mutation scores can be explained by the statement coverage.

Some of the programs in our experiment cannot be explained either by branch or statement coverage. Therefore, there are factors other than the branch coverage and statement coverage that impact mutation score. Those factors are not included in our study and are left for future research to explain program cases such as Jbehave and common-lang. In Jbehave, the rounded mutation score is the highest at 100% although the branch coverage is below 90%. In a contrary example, the branch coverage of common-lang is equal to 90% while its mutation score is 75%. In future work, we will analyze such programs and their associated tests further and incorporate a larger range of TDD and non-TDD developed applications.

9. ACKNOWLEDGMENTS

We would like to thank Dr. René Just, the developer of MAJOR, which we used as our mutation analysis tool, for giving us access to the tool and for his help and guidance in its usage. We would also like to thank Jeshua Kracht and Jacob Petrovic for their assistance in tool usage and debugging. Finally, thanks to Dr. Gregory M. Kapfhammer for his guidance and help in the tools and research for this project.

This should be capitalized.

10. REFERENCES

- [1] Sf100 benchmarks. <http://www.evosuite.org/sf100/>.
- [2] Iec 61508: Functional safety of electrical/electronic/programmable electronic safety-related systems. IEC, 1998.
- [3] James H Andrews, Lionel C Briand, and Yvan Labiche. Is mutation an appropriate tool for testing experiments? [software testing]. In *Software Engineering, 2005. ICSE 2005. Proceedings. 27th International Conference on*, pages 402–411. IEEE, 2005.
- [4] J.H. Andrews, L.C. Briand, Y. Labiche, and A.S. Namin. Using mutation analysis for assessing and comparing testing coverage criteria. *Software Engineering, IEEE Transactions on*, 32(8):608–624, Aug 2006.
- [5] Dave Astels. *Test driven development: A practical guide*. Prentice Hall Professional Technical Reference, 2003.
- [6] Dave Astels. Google techtalk: Beyond test driven development: Behavior driven development, 2006.
- [7] Kent Beck and Cynthia Andres. *Extreme programming explained: embrace change*. Addison-Wesley Professional, 2004.
- [8] Timothy Alan Budd. Mutation analysis of program test data. 1980. ↗ Reference details? words if have 13 room)
- [9] Adnan Cauvic, Sasikumar Punnekkat, and Daniel Sundmark. Quality of testing in test driven development. In *Quality of Information and Communications Technology (QUATIC), 2012 Eighth International Conference on*, pages 266–271. IEEE, 2012.
- [10] Henry Coles. Pit. <http://pitest.org/>, 2014.
- [11] Steve Cornett. Code coverage analysis, 1996–2011.
- [12] Richard A DeMillo, Richard J Lipton, and Frederick G Sayward. Hints on test data selection: Help for the practicing programmer. *Computer*, 11(4):34–41, 1978.
- [13] Michael Feathers. *Working effectively with legacy code*. Prentice Hall Professional, 2004.
- [14] Martin Fowler. *Refactoring: improving the design of existing code*. Addison-Wesley Professional, 1999.
- [15] G. Fraser and A. Zeller. Mutation-driven generation of unit tests and oracles. *Software Engineering, IEEE Transactions on*, 38(2):278–292, March 2012.
- [16] Dick Hamlet and Joe Maybee. *The Engineering of Software: A Technical Guide for the Individual*. Addison-Wesley Longman Publishing Co., Inc., 2000.
- [17] Richard G. Hamlet. Testing programs with the aid of a compiler. *Software Engineering, IEEE Transactions on*, (4):279–290, 1977.
- [18] Marc R Hoffmann and Gilles Iachelini. Code coverage analysis for eclipse. *Eclipse Summit Europe*, 2007.
- [19] René Just, Gregory M Kapfhammer, and Franz Schweiggert. Using conditional mutation to increase the efficiency of mutation analysis. In *Proceedings of the 6th International Workshop on Automation of Software Test*, pages 50–56. ACM, 2011.
- [20] Rene Just, Franz Schweiggert, and Gregory M Kapfhammer. Major: An efficient and extensible tool for mutation analysis in a java compiler. In ↗ Reorder

Every letter in this word should be capitalized.

Proceedings of the ~~2011~~ 26th IEEE/ACM International Conference on Automated Software Engineering, pages 612–615. IEEE Computer Society, 2011.

- [21] Gregory M. Kapfhammer. *The Computer Science and Engineering Handbook*, chapter Chapter 105: Software Testing. CRC Press, Boca Raton, FL, second edition, 2004.
- [22] Clemens Lee. Javancss-a source measurement suite for java, 2005. URL: <http://javancss.codehaus.org>.
- [23] Mark Levinson. Misconceptions with test driven development. <http://agilepainrelief.com/notesfromatooluser/2008/11/misconceptions-with-test-driven-development.html>, Nov 2008.
- [24] Lech Madeyski. The impact of test-first programming on branch coverage and mutation score indicator of unit tests: An experiment. *Information and Software Technology*, 52(2):169–184, 2010.
- [25] Thomas J McCabe. A complexity measure. *Software Engineering, IEEE Transactions on*, (4):308–320, 1976. **Reorder words**
- [26] Co. KG Mountainminds GmbH and Contributors. Jacoco Java code coverage library. <http://www.eclemma.org/>, 2014.
- [27] Nachiappan Nagappan, E Michael Maximilien, Thirumalesh Bhat, and Laurie Williams. Realizing quality improvement through test driven development: results and experiences of four industrial teams. *Empirical Software Engineering*, 13(3):289–302, 2008.
- [28] Dan North. Introducing bdd. *Better Software*, March, 2006.
- [29] A Jefferson Offutt. A practical system for mutation testing: help for the common programmer. In *Test Conference, 1994. Proceedings, International*, pages 824–830. IEEE, 1994. **Reorder words**
- [30] Matjaž Pančur and Mojca Ciglarič. Impact of test-driven development on productivity, code and tests: A controlled experiment. *Information and Software Technology*, 53(6):557–573, 2011.
- [31] SC-167. Do-178b: Software considerations in airborne systems and equipment certification. RTCA, 1992.
- [32] SC-180. Do-254: Design assurance guidance for airborne electronic hardware. RTCA, 2004.
- [33] Elaine J. Weyuker. The evaluation of program-based software test data adequacy criteria. *Communications of the ACM*, 31(6):668–675, 1988.
- [34] Hong Zhu, Patrick AV Hall, and John HR May. Software unit test coverage and adequacy. *AcM computing surveys (CS)*, 29(4):366–427, 1997.

we probably need more comparisons to papers that report on empirical studies

IS this the correct capitalization for these names?

Note: I suggesting cutting certain words from the references (e.g., Maths & publishers) because they are not absolutely needed.

Frank, Weyuker, others...?

This is not the correct capitalization for this tool name.

Ideas for Improving the Paper

1. would it be possible to make the data sets, statistical analysis, graphs, etc. available in a git repository or on a web site?

This would improve the paper's contrib.

and make it possible for others to

I think confirm our results that and replicate our this should experiments. be capitalized.

Should be capitalized

Cutting flex will save space, which is always important in conf & workshop papers!