

Empirical Studies of Test Case Prioritization in a JUnit Testing Environment

Hyunsook Do, Gregg Rothermel, Alex Kinnear
Computer Science and Engineering Department
University of Nebraska - Lincoln
{dohy,grother,akinnear}@cse.unl.edu

Abstract

Test case prioritization provides a way to run test cases with the highest priority earliest. Numerous empirical studies have shown that prioritization can improve a test suite's rate of fault detection, but the extent to which these results generalize is an open question because the studies have all focused on a single procedural language, C, and a few specific types of test suites. In particular, Java and the JUnit testing framework are being used extensively in practice, and the effectiveness of prioritization techniques on Java systems tested under JUnit has not been investigated. We have therefore designed and performed a controlled experiment examining whether test case prioritization can be effective on Java programs tested under JUnit, and comparing the results to those achieved in earlier studies. Our analyses show that test case prioritization can significantly improve the rate of fault detection of JUnit test suites, but also reveal differences with respect to previous studies that can be related to the language and testing paradigm.

1 Introduction

As a software system evolves, software engineers regression test it to detect whether new faults have been introduced into previously tested code. The simplest regression testing technique is to re-run all existing test cases, but this can require a lot of effort, depending on the size and complexity of the system under test. For this reason, researchers have studied various techniques for improving the cost-effectiveness of regression testing, such as regression test selection [10, 27], test suite minimization [9, 18, 23], and test case prioritization [15, 29, 33].

Test case prioritization provides a way to run test cases that have the highest priority — according to some criterion — earliest, and can yield meaningful benefits, such as providing earlier feedback to testers and earlier detection of faults. Numerous prioritization techniques have been described in the research literature and they have been evaluated through various empirical studies [11, 14, 15, 28, 29,

31, 33]. These studies have shown that several prioritization techniques can improve a test suite's rate of fault detection. Most of these studies, however, have focused on a single procedural language, C, and on a few specific types of test suites, so whether their results generalize to other programming and testing paradigms is an open question. Replication of these studies with populations other than those previously examined is needed, to provide a more complete understanding of test case prioritization.

In this work, we set out to perform a replicated study, focusing on an object-oriented language, Java, that is rapidly gaining usage in the software industry, and thus is practically important in its own right. We focus further on a new testing paradigm, the JUnit testing framework, which is increasingly being used by developers to implement test cases for Java programs [1]. In fact, with the introduction of the JUnit testing framework, many software development organizations are building JUnit test cases into their code bases, as is evident through the examination of Open Source Software hosts such as SourceForge and Apache Jakarta [2, 3].

The JUnit framework encourages developers to write test cases, and then to rerun all of these test cases whenever they modify their code. As mentioned previously, however, as the size of a system grows, retest-all regression testing strategies can be excessively expensive. JUnit users will need methodologies with which to remedy this problem.¹

We have therefore designed and performed a controlled experiment examining whether test case prioritization can be effective on object-oriented systems, specifically those written in Java and tested with JUnit test cases. We examine prioritization effectiveness in terms of rate of fault detection, and we also consider whether empirical results show similarity (or dissimilarity) with respect to the results of previous studies. As objects of study we consider four

¹This may not be the case with respect to extreme programming, another development methodology making use of JUnit test cases, because extreme programming is intended for development of modest size projects using a small number of programmers [32]. However, JUnit is also being used extensively in the testing of Java systems constructed by more traditional methodologies, resulting in large banks of integration and system tests. It is in this context that prioritization can potentially be useful.

open source Java programs that have JUnit test suites, and we examine the ability of several test case prioritization techniques to improve the rate of fault detection of these test suites, while also varying other factors that affect prioritization effectiveness. Our results indicate that test case prioritization can significantly improve the rate of fault detection of JUnit test suites, but also reveal differences with respect to previous studies that can be related to the Java and JUnit paradigms.

In the next section of this paper, we describe the test case prioritization problem and related work. Section 3 describes the JUnit testing framework, and our extensions to that framework that allow it to support prioritization. Section 4 presents our experiment design, results, and analysis, describing what we have done in terms of experiment setup to manipulate JUnit test cases. Section 5 discusses our results, and Section 6 presents conclusions and future work.

2 Background and Related Work

2.1 Test Case Prioritization

Test case prioritization techniques [15, 29, 33] schedule test cases in an execution order according to some criterion. The purpose of this prioritization is to increase the likelihood that if the test cases are used for regression testing in the given order, they will more closely meet some objective than they would if they were executed in some other order. For example, testers might schedule test cases in an order that achieves code coverage at the fastest rate possible, exercises features in order of expected frequency of use, or increases the likelihood of detecting faults early in testing.

Depending on the types of information available for programs and test cases, and the way in which those types of information are used, various test case prioritization techniques can be employed. One way in which techniques can be distinguished involves the type of code coverage information they use. Test cases can be prioritized in terms of the number of statements, basic blocks, or methods they executed on a previous version of the software. For example, a *total block coverage prioritization* technique simply sorts test cases in the order of the number of basic blocks (single-entry, single-exit sequences of statements) they covered, resolving ties randomly.

A second way in which prioritization techniques can be distinguished involves the use of “feedback”. When prioritizing test cases, if a particular test case has been selected as “next best”, information about that test case can be used to re-evaluate the value of test cases not yet chosen prior to picking the next test case. For example, *additional block coverage prioritization* iteratively selects a test case that yields the greatest block coverage, then adjusts the coverage information for the remaining test cases to indicate coverage of blocks not yet covered, and repeats this process until

all blocks coverable by at least one test case have been covered. This process is then repeated on remaining test cases.

A third way in which prioritization techniques can be distinguished involves their use of information about code modifications. For example, the amount of change in a code element can be factored into prioritization by weighting the elements covered using a measure of change.

Other dimensions along which prioritization techniques can be distinguished that have been suggested in the literature [14, 15, 22] include test cost estimates, fault severity estimates, estimates of fault propagation probability, test history information, and usage statistics obtained through operational profiles.

2.2 Previous Empirical Work

Early studies of test case prioritization focused on the cost-effectiveness of individual techniques, the estimation of a technique’s performance, or comparisons of techniques [15, 28, 29, 31, 33]. These studies showed that various techniques could be cost-effective, and suggested tradeoffs among them. However, the studies also revealed wide variances in performance, and attributed these to factors involving the programs under test, test suites used to test them, and types of modifications made to the programs.

Recent studies of prioritization have begun to examine the factors affecting prioritization effectiveness [11, 22, 25]. Rothermel et al. [25] studied the effects of test suite design on regression testing techniques, varying the composition of test suites and examining the effects on cost-effectiveness of test selection and prioritization. While this study did not consider correlating attributes of change with technique performance, Elbaum et al. [11] performed experiments exploring characteristics of program structure, test suite composition, and changes on prioritization, and identified several metrics characterizing these attributes that correlate with prioritization effectiveness.

More recent studies have examined how some of these factors affect the effectiveness and efficiency of prioritization, and have considered the generalization of findings through controlled experiments [12, 16, 26]. These studies expose tradeoffs and constraints that affect the success of techniques, and provide guidelines for designing and managing prioritization and testing processes.

Most recently, Saff and Ernst [30] considered test case prioritization for Java in the context of continuous testing, which uses spare CPU resources to continuously run regression tests in the background as a programmer codes. They propose combining the concepts of test frequency and test case prioritization, and report the results of a study that show that prioritized continuous testing reduced wasted development time. However, their prioritization techniques are based on different sources of information than ours, such as history of recent or frequent errors and test cost,

rather than code coverage information. The measure of effectiveness they use also differs from ours: theirs involves reduction of wasted time in development, whereas ours involves the weighted average of the percentage of faults detected over the life of a test suite.

With the exception of the work reported in [30], all of this previous empirical work has concerned C programs and system-level test suites constructed for code coverage, or for partition-based coverage of requirements. In contrast, the study we describe here examines whether prior results generalize, by replicating previous experiments on a new population of programs and test suites (Java and JUnit), and examining whether the results are consistent with those of the previous studies.

3 JUnit Testing and Prioritization

JUnit test cases are Java classes that contain one or more test methods and that are grouped into test suites, as shown in Figure 1. The figure presents a simple hierarchy having only a single test-class level, but the tree can extend deeper using additional nesting of Test Suites. The leaf nodes in such a hierarchy, however, always consist of test-methods, where a test-method is a minimal unit of test code.

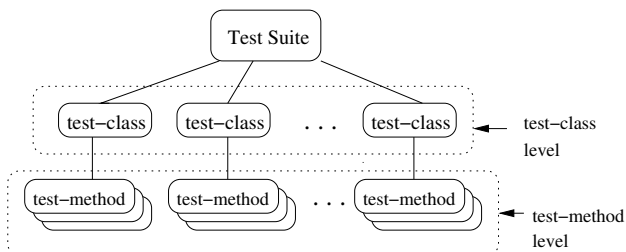


Figure 1. JUnit test suite structure

JUnit test classes that contain one or more test methods can be run individually, or a collection of JUnit test cases (a test suite) can be run as a unit. Running individual test cases is reasonable for small programs, but for large numbers of test cases can be expensive, because each independent execution of a test case incurs startup costs. Thus in practice, developers design JUnit test cases to run as sequences of tests invoked through test suites that invoke test classes.

Clearly, choices in *test suite granularity* (the number and size of the test cases making up a test suite) can affect the cost of running JUnit test cases, and we want to investigate the relationship between this factor and prioritization results. To do this, we focus on two levels of test suite granularity: test-class level, a collection of test-classes that represents a coarse test suite granularity, and test-method level, a collection of test-methods that represents a fine test suite granularity. To support this focus we needed to ensure that the JUnit framework allowed us to achieve the following four objectives:

1. treat each TestCase class as a single test case for purposes of prioritization (test-class level);
2. reorder TestCase classes to produce a prioritized order;
3. treat individual test methods within TestCase classes as test cases for prioritization (test-method level);
4. reorder test methods to produce a prioritized order.

Objectives 1 and 2 were trivially achieved as a consequence of the fact that the default unit of test code that can be specified for execution in the JUnit framework is a TestCase class. Thus it was necessary only to extract the names of all TestCase classes invoked by the top level TestSuite for the object program² (a simple task) and then execute them individually with the JUnit test runner in a desired order.

Objectives 3 and 4 were more difficult to achieve, due to the fact that a TestCase class is also the minimal unit of test code that can be specified for execution in the normal JUnit framework. Since a TestCase class can define multiple test methods, all of which will be executed when the TestCase is specified for execution, providing the ability to treat individual methods as test cases required us to extend the JUnit framework to support this finer granularity. Thus the principal challenge we faced was to design and implement JUnit extensions that provide a means for specifying individual test methods for execution, as found in the total set of methods distributed across multiple TestCase classes.

Since the fundamental purpose of the JUnit framework is to discover and execute test methods defined in TestCase classes, the problem of providing test-method level testing reduces to the problem of uniquely identifying each test method discovered by the framework and making them available for individual execution by the tester. We accomplished this task by extending (subclassing) various components of the framework and inserting mechanisms for assigning numeric test IDs to each test method discovered. We then created a SelectiveTestRunner that uses the new extension components. The relationship between our extensions and the existing JUnit framework is shown in Figure 2, which also shows how the JUnit framework is related to the Galileo system for analyzing Java bytecode (which we used to obtain coverage information for use in prioritization). Our new SelectiveTestRunner is able to access test cases individually using numeric test IDs.

To implement prioritization at the test-method level we also needed to provide a way for the test methods to be executed in a tester-specified order. Because the JUnit framework must discover the test methods, and our extensions assign numeric IDs to tests in the order of discovery, to execute the test cases in an order other than the one in which they are provided requires that all test cases be discovered prior to execution. We accomplished this by using a simple two-pass technique. In the first pass, all the test methods rel-

²This process may need to be repeated iteratively if Test Suites are nested in other Test Suites.

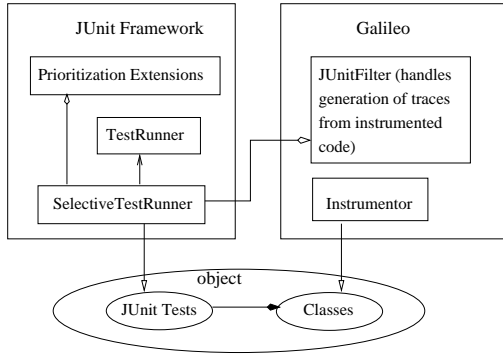


Figure 2. JUnit framework and Galileo

evant to the system are discovered and assigned numbers. The second pass then uses a specified ordering to retrieve and execute each method by its assigned ID. The tester (or testing tool) is provided a means, via the SelectiveTestRunner, of retrieving the IDs assigned to test methods by the framework for use with prioritization.

4 Experiments

We wish to address the following research questions:

- RQ1: Can test case prioritization improve the rate of fault detection of JUnit test suites?
- RQ2: How do the three types of information and information use that distinguish prioritization techniques (type of coverage information, use of feedback, and use of modification information) impact the effectiveness of prioritization techniques?
- RQ3: Can test suite granularity (the choice of running test-class level versus test-method level JUnit test cases) impact the effectiveness of prioritization techniques?

In addition to these research questions, we examine whether test case prioritization results obtained from systems written in an object-oriented language (Java) and using JUnit test cases show different trends than those obtained from systems written in a procedural language (C) using traditional coverage- or requirements-based system test cases.

To address our questions we designed several controlled experiments. The following subsections present, for these experiments, our objects of analysis, independent variables, dependent variables and measures, experiment setup and design, threats to validity, and data and analysis.

4.1 Objects of Analysis

We used four Java programs as objects of analysis: Ant, XML-security, JMeter, and JTopas. Ant is a Java-based build tool [4]; it is similar to make, but instead of being extended with shell-based commands it is extended using

Java classes. JMeter is a Java desktop application designed to load test functional behavior and measure performance [5]. XML-security implements security standards for XML [6]. JTopas is a Java library used for parsing text data [7]. Several sequential versions of each of these systems were available and were selected for these experiments.

Table 1 lists, for each of our objects, “Versions” (the number of versions), “Size” (the number of classes), “Test Classes” (the number of JUnit test-class level test cases), “Test Methods” (the number of JUnit test-method level test cases), and “Faults” (the number of faults). The number of classes corresponds to the total number of class files in the final version. The numbers for test-class (test-method) list the number of test classes (test methods) in the most recent version. The number of faults indicates the total number of faults available for each of the objects (see Section 4.3.3).

Table 1. Experiment objects

Subjects	Versions	Size	Test Classes	Test Methods	Faults
Ant	9	627	150	877	21
XML-security	4	143	14	83	6
JMeter	6	389	28	78	9
JTopas	4	50	11	128	5

4.2 Variables and Measures

4.2.1 Independent Variables

Our experiments manipulated two independent variables: prioritization technique and test suite granularity.

Variable 1: Prioritization Technique

We consider nine different test case prioritization techniques, which we classify into three groups to match an earlier study on prioritization for C programs [15]. Table 2 summarizes these groups and techniques. The first group is the control group, containing three “techniques” that serve as experimental controls. (We use the term “technique” here as a convenience; in actuality, the control group does not involve any practical prioritization heuristics; rather, it involves various orderings against which practical heuristics should be compared.) The second group is the block level group, containing two fine granularity prioritization techniques. The third group is the method level group, containing four coarse granularity prioritization techniques.

Control techniques

- No prioritization (untreated): One control that we consider is simply the application of no technique; this lets us consider “untreated” JUnit test suites.
- Random prioritization (random): As a second control we use random prioritization, in which we randomly order the test cases in a JUnit test suite.

Table 2. Test case prioritization techniques.

Label	Mnemonic	Description
T1	untreated	original ordering
T2	random	random ordering
T3	optimal	ordered to optimize rate of fault detection
T4	block-total	prioritize on coverage of block
T5	block-addtl	prioritize on coverage of block not yet covered
T6	method-total	prioritize on coverage of method
T7	method-addtl	prioritize on coverage of method not yet covered
T8	method-diff-total	prioritize on coverage of method and change information
T9	method-diff-addtl	prioritize on coverage of method/change information, and adjusted on previous coverage

- Optimal prioritization (optimal): To measure the effects of prioritization techniques on rate of fault detection, our empirical study uses programs that contain known faults. For the purposes of experimentation we can determine, for any test suite, which test cases expose which faults, and thus we can determine an optimal ordering of test cases in a JUnit test suite for maximizing that suite’s rate of fault detection. This is not a viable practical technique, but it provides an upper bound on the effectiveness of our heuristics.

Block level techniques

- Total block coverage prioritization (block-total): By instrumenting a program we can determine, for any test case, the number of basic blocks in that program that are exercised by that test case. We can prioritize these test cases according to the total number of blocks they cover simply by sorting them in terms of that number.
- Additional block coverage prioritization (block-addtl): Additional block coverage prioritization combines *feedback* with coverage information. It iteratively selects a test case that yields the greatest block coverage, adjusts the coverage information on subsequent test cases to indicate their coverage of blocks not yet covered, and repeats this process until all blocks covered by at least one test case have been covered. If multiple test cases cover the same number of blocks not yet covered, they are ordered randomly. When all blocks have been covered, this process is repeated on the remaining test cases until all have been ordered.

Method level techniques

- Total method coverage prioritization (method-total): Total method coverage prioritization is the same as total block coverage prioritization, except that it relies on coverage measured in terms of methods.

- Additional method coverage prioritization (method-addtl): Additional method coverage prioritization is the same as additional block coverage prioritization, except that it relies on coverage in terms of methods.
- Total diff method prioritization (method-diff-total): Total diff method coverage prioritization uses modification information; it sorts test cases in the order of their coverage of methods that differ textually (as measured by a Java parser that parses pairs of individual Java methods through the Unix “diff” function). If multiple test cases cover the same number of differing methods, they are ordered randomly.
- Additional diff method prioritization (method-diff-addtl): Additional diff method prioritization uses both feedback and modification information. It iteratively selects a test case that yields the greatest coverage of methods that differ, adjusts the information on subsequent test cases to indicate their coverage of methods not yet covered, and then repeats this process until all methods that differ and have been covered by at least one test case have been covered. If multiple test cases cover the same number of differing methods not yet covered, they are ordered randomly. This process is repeated until all test cases that execute methods that differ have been used; additional method coverage prioritization is applied to remaining test cases.

The foregoing set of techniques matches the set examined in [15] in all but two respects. First, we use three control techniques, considering an “untreated” technique in which test cases are run in the order in which they are given in the original JUnit test cases. This is a sensible control technique for our study since in practice developers would run JUnit test cases in their original ordering.

Second, the studies with C programs used statement and function level prioritization techniques, where coverage is based on source code, whereas our study uses coverage based on Java bytecode. Analysis at the bytecode level is appropriate for Java environments. Since Java is a platform independent language, vendors or programmers might choose to provide just class files for system components. In such cases we want to be able to analyze even those class files, and bytecode analysis allows this.

The use of bytecode level analysis does affect our choice of prioritization techniques. As an equivalent to C “function level” coverage, a method level granularity was an obvious choice. As a statement level equivalent, we could use either individual bytecode instructions, or basic blocks of instructions, but we cannot infer a one-to-one correspondence between Java source statements and either bytecode instructions or blocks.³ We chose the basic block because the basic

³A Java source statement typically compiles to several bytecode instructions, and a basic block from bytecode often corresponds to more than one Java source code statement.

block representation is a more cost-effective unit of analysis for bytecode. Since basic blocks of bytecode and source code statements represent different levels of granularity, results obtained on the two are not directly comparable, but we can still study the effect that an increase in granularity has on prioritization in Java, as compared to in C.

Variable 2: Test Suite Granularity

To investigate the impact of test suite granularity on the effectiveness of test case prioritization techniques we considered two test suite granularity levels for JUnit test cases: test-class level and test-method level as described in Section 3. At the test-class level, JUnit test cases are of relatively coarse granularity; each test suite and test-class that is invoked by a TestSuite is considered to be one test case, consisting of one or more test-methods. At the test-method level, JUnit test cases are relatively fine granularity; each test-method is considered to be one test case.

4.2.2 Dependent Variables and Measures

Rate of Fault Detection

To investigate our research questions we need to measure the benefits of the various prioritization techniques in terms of rate of fault detection. To measure rate of fault detection, we use a metric, APFD (Average Percentage Faults Detected), introduced for this purpose in [15], that measures the weighted average of the percentage of faults detected over the life of a test suite. APFD values range from 0 to 100; higher numbers imply faster (better) fault detection rates. More formally, let T be a test suite containing n test cases, and let F be a set of m faults revealed by T . Let TF_i be the first test case in ordering T' of T which reveals fault i . The APFD for test suite T' is given by the equation:

$$APFD = 1 - \frac{TF_1 + TF_2 + \dots + TF_m}{nm} + \frac{1}{2n}$$

Examples and empirical results illustrating the use of this metric are provided in [15].

4.3 Experiment Setup

To perform test case prioritization we required several types of data. Since the process used to collect this data is complicated and requires significant time and effort, we automated a large part of the experiment.

Figure 3 illustrates our experiment process. There were three types of data to be collected prior to applying prioritization techniques: coverage information, fault-matrices, and change information. We obtained coverage information by running test cases over instrumented objects using the Galileo system for analysis of Java bytecode in conjunction with a special JUnit adaptor. This information lists which

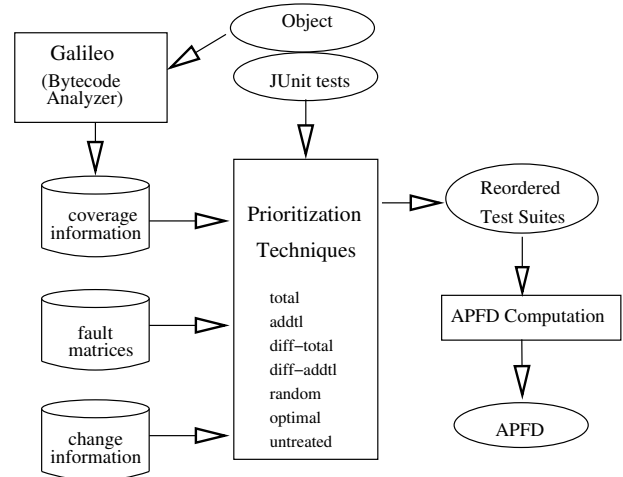


Figure 3. Overview of experiment process

test cases exercised which blocks and methods; a previous version’s coverage information is used to prioritize the current set of test cases. Fault-matrices list which test cases detect which faults and are used to measure the rate of fault detection for each prioritization technique. Change information lists which methods differ from those in the preceding version and how many lines of each method were changed (deleted and added methods are also listed.).

Each prioritization technique uses some or all of this data to prioritize JUnit test suites based on its analysis; then APFD scores are obtained from the reordered test suites. The collected scores are analyzed to determine whether the techniques improved the rate of fault detection.

4.3.1 Object Instrumentation

To perform our experiment, we required objects to be instrumented to support the techniques described in Section 4.2.1. We instrumented object class files in two ways: all basic blocks, and all method entry blocks, using the Galileo bytecode analysis system (see Figure 2).

4.3.2 Test Cases and Test Automation

As described previously, test cases were obtained from each object’s distribution, and we considered test suites at the two levels of granularity previously described. To execute and validate test cases automatically, we created test scripts that invoke test cases, save all outputs from test execution, and compare outputs with those for the previous version. As shown in Figure 2, JUnit test cases are run through JUnitFilter and TestRunner (SelectiveTestRunner) over the instrumented classes.

4.3.3 Faults

We wished to evaluate the performance of prioritization techniques with respect to detection of regression faults. The object programs we obtained were not supplied with any such faults or fault data. To obtain faults, we considered two approaches: mutation and fault seeding. The first approach would allow us to generate a large number of faults, but these faults would not be representative of real faults. The second approach cannot cost-effectively produce a large number of faults, but it can generate more realistic faults than those created by mutation. Thus, we chose the second approach, and following the procedure described in [20], and also used in the study described in [15], we seeded faults. Two graduate students performed this fault seeding; they were instructed to insert faults that were as realistic as possible based on their experience with real programs, and that involved code inserted into, or modified in each of the versions. To be consistent with previous studies, we excluded any faults that were detected by more than 20% of the test cases at both granularity levels.

4.4 Threats to Validity

In this section we describe the internal, external, and construct threats to the validity of our experiments, and the approaches we used to limit the effects that these threats might have.

Internal Validity

The inferences we have made about the effectiveness of prioritization techniques could have been affected by two factors. The first factor involves the potential faults in our experiment tools. To control for this threat, we validated tools through testing on various sizes of Java programs. The second factor involves the faults seeded in our objects. Our procedure for seeding faults followed a set process as described in Section 4.3.3, which reduced the chances of obtaining biased faults. However some of our objects (XML-security and JTopas) ultimately contained a relatively small number of faults, and this might affect results.

External Validity

Three issues limit the generalization of our results. The first issue is object program representativeness. Our objects are of small and medium size. Complex industrial programs with different characteristics may be subject to different cost-benefit tradeoffs. The second issue involves testing process representativeness. If the testing process we used is not representative of industrial processes, our results might not generalize. Control for these threats can be achieved only through additional studies with a wider population. The third issue involves fault representativeness. We used seeded faults that were as realistic as possible, but they

were not real faults and were limited to one type of artificial fault. Future study will need to consider other fault types.

Construct Validity

The dependent measure that we have considered, APFD, is not the only possible measure of rate of fault detection and has some limitations. For example, APFD assigns no value to subsequent test cases that detect a fault already detected; such inputs may, however, help debuggers isolate the fault, and for that reason might worth measuring. Also, APFD does not account for the possibility that faults and test cases may have different costs. Future studies will need to consider other measures for purposes of assessing effectiveness.

4.5 Data and Analysis

To provide an overview of all the collected data we present boxplots in Figure 4. The left side of the figure presents results from test case prioritization applied to test-class level test cases, and the right side presents results from test case prioritization applied to test-method level test cases. The top row presents results for an all programs total, and other rows present results per object program. Each plot contains a box showing the distribution of APFD scores for each of the nine techniques, across each of the versions of the object program. See Table 2 for a legend of the techniques.

The data sets depicted in Figure 4 served as the basis for our formal analyses of results. The following sections describe, for each of our research questions in turn, the experiments relevant to that question, presenting those analyses.

4.5.1 RQ1: Prioritization effectiveness

Our first research question considers whether test case prioritization can improve the rate of fault detection for JUnit test cases applied to our Java objects.

An initial indication of how each prioritization technique affected a JUnit test suite's rate of fault detection in this study can be obtained from Figure 4. Comparing the boxplots of optimal (T3) to those for untreated (T1) and random (T2), it is apparent that an optimal prioritization order could provide substantial improvements in rates of fault detection. Comparing results for untreated (T1) to results for actual, non-control techniques (T4 through T9) for both test suite levels, it appears that all non-control techniques yield improvement. However, the comparison of the results of random orderings (T2) with those produced by non-control techniques shows different results: some techniques yield improvement with respect to random while others do not.

To determine whether the differences observed in the boxplots are statistically significant we performed two sets

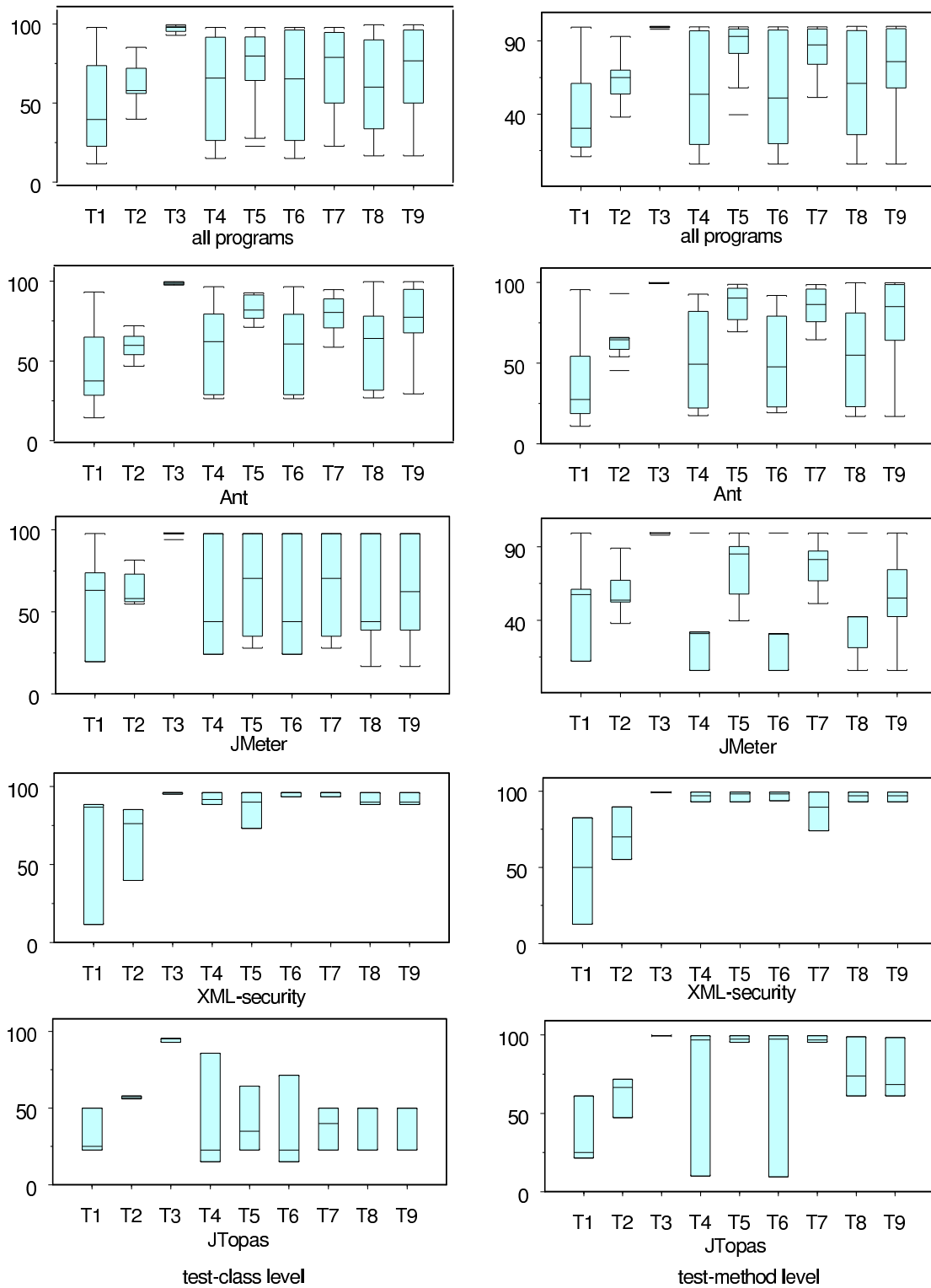


Figure 4. APFD boxplots, all programs. The horizontal axes list techniques, and the vertical axes list APFD scores. The left column presents results for test-class level test cases and the right column presents results for test-method level test cases.

test-class level					
Source	SS	d.f	MS	F-stat	p-value
Technique	7943	6	1323	2.09	0.058
Program	25942	3	8647	13.7	< 0.0001
Residuals	77584	123	630		
multiple comparison with a control by Dunnett's method					
critical point: 2.603					
	Estimate	Std. Err	Lower Bound	Upper Bound	
T4-T1	11.7	8.15	-9.47	32.9	
T5-T1	23.5	8.15	2.28	44.7	****
T6-T1	11.3	8.15	-9.91	32.5	
T7-T1	22.8	8.15	1.56	44.0	****
T8-T1	12.6	8.15	-8.65	33.8	
T9-T1	20.6	8.15	-0.61	41.8	
test-method level					
Source	SS	d.f	MS	F-stat	p-value
Technique	30711	6	5118	7.02	< 0.0001
Program	18902	3	6300	8.64	< 0.0001
Residuals	89637	123	728		
multiple comparison with a control by Dunnett's method					
critical point: 2.603					
	Estimate	Std. Err	Lower Bound	Upper Bound	
T4-T1	15.0	8.76	-7.83	37.8	
T5-T1	44.6	8.76	21.80	67.4	****
T6-T1	14.6	8.76	-8.16	37.4	
T7-T1	43.0	8.76	20.20	65.8	****
T8-T1	19.3	8.76	-3.49	42.1	
T9-T1	31.8	8.76	9.02	54.6	****

Table 3. ANOVA and Dunnett's method, untreated vs heuristics, all programs.

test-class level					
Source	SS	d.f	MS	F-stat	p-value
Technique	3481	6	492	1.06	0.389
Program	23262	3	7754	14.1	< 0.0001
Residuals	67272	123	546		
multiple comparison with a control by Dunnett's method					
critical point: 2.603					
	Estimate	Std. Err	Lower Bound	Upper Bound	
T4-T2	-1.140	7.59	-20.90	18.6	
T5-T2	10.600	7.59	-9.14	30.4	
T6-T2	-1.580	7.59	-21.30	18.2	
T7-T2	9.880	7.59	-9.87	29.6	
T8-T2	-0.326	7.59	-20.10	19.4	
T9-T2	7.720	7.59	-12.00	27.5	
test-method level					
Source	SS	d.f	MS	F-stat	p-value
Technique	18519	6	3086	4.94	< 0.0001
Program	20613	3	6871	10.99	< 0.0001
Residuals	76848	123	624		
multiple comparison with a control by Dunnett's method					
critical point: 2.603					
	Estimate	Std. Err	Lower Bound	Upper Bound	
T4-T2	-6.85	8.11	-28.0	14.3	
T5-T2	22.80	8.11	1.72	43.9	****
T6-T2	-7.18	8.11	-28.3	13.9	
T7-T2	21.20	8.11	0.09	42.3	****
T8-T2	-2.50	8.11	-23.6	18.6	
T9-T2	10.00	8.11	-11.1	31.1	

Table 4. ANOVA and Dunnett's method, random vs heuristics, all programs.

of analyses, considering test-class and test-method levels independently for all programs. The analyses were:

1. **UNTREATED vs NON-CONTROL:** We consider untreated and non-control techniques to determine whether there is a difference between techniques (using ANOVA), and whether there is a difference between untreated orders and the orders produced by each non-control technique (using Dunnett's method for multiple comparison with a control [19]).
2. **RANDOM vs NON-CONTROL:** We perform the same analyses as in (1), against the random ordering.

Table 3 presents the results of analysis (1), for a significance level of 0.05. The ANOVA results indicate that there is no difference between the techniques at the test-class level (p-value = 0.058), but there is a significant difference at the test-method level (p-value < 0.0001). Dunnett's method reports that block-addtl and method-addtl are different from untreated at the test-class level, and block-addtl, method-addtl, and method-diff-addtl are different from untreated at the test-method level. Cases that were statistically

significant are marked by "****" (which indicates confidence intervals that do not include zero).

Table 4 presents the results of analysis (2). Similar to the first analysis, the ANOVA results indicate that there is no difference between the techniques at the test-class level (p-value = 0.389), but there is a significant difference at the test-method level (p-value < 0.0001). Dunnett's method reports that no individual prioritization technique differs from random at the test-class level, but block-addtl and method-addtl differ from random at the test-method level.

4.5.2 RQ2: The effects of information types and use on prioritization results

Our second research question concerns whether differences in the types of information and information use that distinguish prioritization techniques (type of coverage information, use of feedback, type of modification information) impact the effectiveness of prioritization.

Comparing the boxplots (Figure 4) of block-total to method-total and block-addtl to method-addtl for both test suite levels, it appears that the level of coverage information utilized (fine vs coarse) has no effect on techniques. Com-

paring the results of block-total to block-addtl and method-total to method-addtl for both test suite levels, it appears that techniques using feedback do yield improvement over those not using feedback. Finally, comparison of the results of method-total to method-diff-total and method-addtl to method-diff-addtl shows no apparent effect from using modification information on prioritization effectiveness.

To determine whether the differences observed in the boxplots are statistically significant we performed an ANOVA and all-pairwise comparison using Tukey’s method [8], considering test-class level and test-method level independently, for all non-control techniques and all programs. Table 5 presents the results of these analyses.

Overall, the results indicate that there is a significant difference between techniques at the test-method level (p-value = 0.00022); but no significant difference between techniques at the test-class level (p-value = 0.377).

The all-pairwise comparison shows the following with respect to information types and information use:

- Coverage information. The results indicate that there is no difference between block-level and method-level techniques at either test suite level, considering block-total versus method-total (T4 - T6) and block-addtl versus method-addtl (T5 - T7). Different levels of coverage information did not impact the effectiveness of prioritization.
- Use of Feedback. The results indicate a significant difference between techniques that use feedback and those that do not use feedback at the test-method level, namely block-total versus block-addtl (T4 - T5) and method-total versus method-addtl (T6 - T7). However, there was no significant difference at the test-class level.
- Modification information. The results indicate no significant difference between techniques that use modification information and those that do not use modification information, namely method-total versus method-diff-total (T6 - T8) and method-addtl versus method-diff-addtl (T7 - T9), at either test suite level.

4.5.3 RQ3: Test suite granularity effects

Our third research question considers the impact of test suite granularity, comparing test-class level to test-method level. The boxplots and the analysis related to our first two research question suggest that there is a difference between the two levels of test suite granularity, thus we performed an ANOVA for non-control techniques and all programs comparing test-method level to test-class level. The results indicate that there is suggestive evidence that test suite granularity affected the rate of fault detection for block-addtl (T5, p-value = 0.038) and method-addtl (T7, p-value = 0.048) techniques.

test-class level					
Source	SS	d.f	MS	F-stat	p-value
Technique	3195	5	639	1.076	0.377
Program	25738	3	8579	14.44	< 0.0001
Residuals	62346	105	593		
all pair comparison by Tukey’s method					
critical point: 2.903					
	Estimate	Std. Err	Lower Bound	Upper Bound	
T4-T5	-11.700	7.91	-34.7	11.2	
T4-T6	0.438	7.91	-22.5	23.4	
T4-T7	-11.000	7.91	-34.0	11.9	
T4-T8	-0.817	7.91	-23.8	22.1	
T4-T9	-8.860	7.91	-31.8	14.1	
T5-T6	12.200	7.91	-10.8	35.1	
T5-T7	0.723	7.91	-22.2	23.7	
T5-T8	10.900	7.91	-12.0	33.9	
T5-T9	2.890	7.91	-20.1	25.8	
T6-T7	-11.500	7.91	-34.4	11.5	
T6-T8	-1.260	7.91	-24.2	21.7	
T6-T9	-9.300	7.91	-32.2	13.7	
T7-T8	10.200	7.91	-12.7	33.2	
T7-T9	2.170	7.91	-20.8	25.1	
T8-T9	-8.040	7.91	-31.0	14.9	
test-method level					
Source	SS	d.f	MS	F-stat	p-value
Technique	17882	5	3576	5.28	0.00022
Program	22367	3	7455	11.0	< 0.0001
Residuals	71112	105	677		
all pair comparison by Tukey’s method					
critical point: 2.903					
	Estimate	Std. Err	Lower Bound	Upper Bound	
T4-T5	-29.700	8.44	-54.000	-5.34	****
T4-T6	0.332	8.44	-24.000	24.70	
T4-T7	-28.100	8.44	-52.400	-3.71	****
T4-T8	-4.340	8.44	-28.700	20.00	
T4-T9	-16.800	8.44	-41.200	7.49	
T5-T6	30.000	8.44	5.670	54.40	****
T5-T7	1.630	8.44	-22.700	26.00	
T5-T8	25.300	8.44	0.998	49.70	****
T5-T9	12.800	8.44	-11.500	37.20	
T6-T7	-28.400	8.44	-52.700	-4.04	****
T6-T8	-4.670	8.44	-29.000	19.70	
T6-T9	-17.200	8.44	-41.500	7.16	
T7-T8	23.700	8.44	-0.631	48.00	
T7-T9	11.200	8.44	-13.100	35.50	
T8-T9	-12.500	8.44	-36.800	11.80	

Table 5. ANOVA and Tukey’s method, all heuristics , all programs.

5 Discussion

Our results strongly support the conclusion that test case prioritization techniques can improve the rate of fault detection for JUnit test suites applied to Java systems. The prioritization techniques we examined outperformed both untreated and randomly ordered test suites, as a whole, at the test-method level. Overall, at the test-class level, prioritization techniques did not improve effectiveness compared

to untreated or randomly ordered test suites, but individual comparisons indicated that techniques using additional coverage information did improve the rate of fault detection.

We also observed that random test case orderings outperformed untreated test case orderings at both test suite levels. We conjecture that this difference is due to the construction of the JUnit test cases supplied with the objects used in this study. It is typical in practice for developers to add new test cases at the end of a test suite. Since newer test cases tend to exercise new code, we conjecture that these new test cases are more likely to be fault-revealing than previous test cases. Randomly ordered test cases redistribute fault-revealing test cases more evenly than original, untreated test case orders. Of course prioritization techniques, in turn, typically outperform random orderings. The practical implication of this result, then, is that the worst thing JUnit users can do is not practice some form of prioritization.

We now consider whether the results we have observed are consistent with those from the previous studies with C programs and coverage- or requirements-based test suites. Both this study (at the test-method level) and previous studies [13, 15, 25, 29] showed that prioritization techniques improved the rate of fault detection compared to both random and untreated orderings. Also, both this study and earlier studies found that techniques using additional coverage information were usually better than other techniques, for both fine and coarse granularity test cases. There were sources of variation between the studies, however.

Regarding the impact of granularity, previous studies on C showed that statement-level techniques as a whole were better than function-level techniques. This study, however, found that block-level techniques were not significantly different overall from method-level techniques. There are two possible reasons for this difference. First, this result may be due to the fact that the instrumentation granularity we used for Java programs differs from that used for C programs, as we explained in Section 4.2.1. Block-level instrumentation is not as sensitive as statement-level instrumentation since a block may combine a sequence of consecutive statements into a single unit in a control flow graph.

A second related factor that might account for this result is that the instrumentation difference between blocks and methods is not as pronounced in Java as is the difference between statements and functions in C. Java methods tend to be more concise than C functions, possibly due to object-oriented language characteristics [21] and code refactoring [17], which tend to result in methods containing small numbers of blocks. For example, constructors, “get” methods, and “set” methods frequently contain only one basic block. A study reported in [24] supports this interpretation, providing evidence that the sizes of the methods called most frequently in object-oriented programs are between one and nine statements on average, which generally corresponds

(in our measurements on our objects) to only one or two basic blocks. Since instrumentation at the method level is less expensive than instrumentation at the basic-block level, if these results generalize, these results have implications for practitioners attempting to select a technique.

Where test suite granularity effects are concerned, contrary to the previous study [25], this study showed that fine granularity test suites (test-method level) were more supportive of prioritization than coarse granularity test suites (test-class level). Although some techniques (block-`addtl` and method-`addtl`) at the test-class level improved rate of fault detection, overall analysis showed that test-class level test suites did not support improvements in APFD. Since the scope of each test case in a test-class level test suite is limited to a specific class under test,⁴ one possible explanation for this result involves the fact that, in our objects, faults are located only in a few of the classes under test, and the number of faults in several cases is relatively small. The test-class level boxplots for the Ant program (Figure 4), which has a larger number of faults (21) overall, show a different trend from the other plots, supporting this notion.

6 Conclusions and Future Work

We have presented a study of prioritization techniques applied across four Java programs, to JUnit test suites provided with those programs. Although several studies of test case prioritization have been conducted previously, most studies have focused on a single procedural language, C, and on only a few specific types of test suites. Our study, in contrast, applied prioritization techniques to an object-oriented language (Java) tested under the JUnit testing framework, to investigate whether the results observed in previous studies generalize to other language and testing paradigms.

Our results regarding the effectiveness of prioritization techniques confirm several previous findings [13, 15, 25, 29], while also revealing some differences regarding prioritization technique granularity effects and test suite granularity effects. As discussed in Section 5, these differences can be explained in relation to characteristics of the Java language and JUnit testing.

The results of our studies suggest several avenues for future work. First, we intend to perform additional studies using larger Java programs and additional types of test suites, and a wider range and distribution of faults. Second, since our studies, unlike previous studies, did not find differences between prioritization technique granularity levels, it would be interesting to investigate what types of attributes caused this result through a further analysis of the relationship between basic block size and method size, and the ratio of execution of simple methods to complex methods. Third, con-

⁴This is true in the objects that we considered; in general, however, the scope of a unit test could be subject to a developers’ specific practices.

trary to our expectation, code modification information did not improve cost-effectiveness significantly. This result was also observed in C programs, and it is worth investigating what types of factors other than fault distribution in a code base could be involved in this outcome. Fourth, because our analysis revealed a sizable performance gap between prioritization heuristics and optimal prioritization, we are investigating alternative techniques. One such alternative involves employing measures of static and dynamic dependencies in the code to estimate locations where faults can reside.

Through the results reported in this paper, and this future work, we hope to provide software practitioners, in particular practitioners who use Java and the JUnit testing framework, with cost-effective techniques for improving regression testing processes through test case prioritization.

Acknowledgements

This work was supported in part by NSF under Awards, CCR-0080900 and CCR-0306023 to Oregon State University, where much of this work was performed. Ziru Zhu helped prepare the JTopas object.

References

- [1] <http://www.junit.org>.
- [2] <http://sourceforge.net>.
- [3] <http://jakarta.apache.org>.
- [4] <http://ant.apache.org>.
- [5] <http://jakarta.apache.org/jmeter>.
- [6] <http://xml.apache.org/security>.
- [7] <http://jtopas.sourceforge.net/jtopas>.
- [8] Splus statistics software.
- [9] T. Chen and M. Lau. Dividing strategies for the optimization of a test suite. *Information Processing Letters*, 60(3):135–141, Mar. 1996.
- [10] Y. Chen, D. Rosenblum, and K. Vo. TestTube: A system for selective regression testing. In *Proc. Int'l. Conf. Softw. Eng.*, pages 211–220, May 1994.
- [11] S. Elbaum, D. Gable, and G. Rothermel. Understanding and measuring the sources of variation in the prioritization of regression test suites. In *Proc. Int'l. Softw. Metrics Symp.*, pages 169–179, Apr. 2001.
- [12] S. Elbaum, P. Kallakuri, A. Malishevsky, G. Rothermel, and S. Kanduri. Understanding the effects of changes on the cost-effectiveness of regression testing techniques. *J. Softw. Testing, Verif., and Rel.*, 12(2), 2003.
- [13] S. Elbaum, A. Malishevsky, and G. Rothermel. Prioritizing test cases for regression testing. In *Proc. Int'l. Symp. Softw. Testing and Anal.*, pages 102–112, Aug. 2000.
- [14] S. Elbaum, A. Malishevsky, and G. Rothermel. Incorporating varying test costs and fault severities into test case prioritization. In *Proc. Int'l. Conf. Softw. Eng.*, pages 329–338, May 2001.
- [15] S. Elbaum, A. G. Malishevsky, and G. Rothermel. Test case prioritization: A family of empirical studies. *IEEE Trans. Softw. Eng.*, 28(2):159–182, Feb. 2002.
- [16] S. Elbaum, G. Rothermel, S. Kanduri, and A. G. Malishevsky. Selecting a cost-effective test case prioritization technique. Technical Report 03-01-01, University of Nebraska – Lincoln, Jan. 2003.
- [17] M. Fowler. *Refactoring: Improving the Design of Existing Code*. Addison-Wesley, 1999.
- [18] M. J. Harrold, R. Gupta, and M. L. Soffa. A methodology for controlling the size of a test suite. *ACM Trans. Softw. Eng. and Meth.*, 2(3):270–285, July 1993.
- [19] J. Hsu. *Multiple Comparisons: Theory and Methods*. Chapman & Hall, London, 1996.
- [20] M. Hutchins, H. Foster, T. Goradia, and T. Ostrand. Experiments on the effectiveness of dataflow- and controlflow-based test adequacy criteria. In *Proc. Int'l. Conf. Softw. Eng.*, pages 191–200, May 1994.
- [21] K. Ishizaki, M. Kawahito, T. Yasue, M. Takeuchi, T. Ogasawara, T. Suganuma, T. Onodera, H. Komatsu, and T. Nakatani. Design, implementation and evaluation of optimisations in a just-in-time compiler. In *ACM 1999 Java Grande Conf.*, pages 119–128, June 1999.
- [22] J. Kim and A. Porter. A history-based test prioritization technique for regression testing in resource constrained environments. In *Proc. Int'l. Conf. Softw. Eng.*, May 2002.
- [23] J. Offutt, J. Pan, and J. M. Voas. Procedures for reducing the size of coverage-based test sets. In *Proc. Int'l. Conf. Testing Comp. Softw.*, pages 111–123, June 1995.
- [24] J. Power and J. Waldron. A method-level analysis of object-oriented techniques in Java applications. Technical Report NUM-CS-TR-2002-07, National University of Ireland, July 2002.
- [25] G. Rothermel, S. Elbaum, A. Malishevsky, P. Kallakuri, and B. Davia. The impact of test suite granularity on the cost-effectiveness of regression testing. In *Proc. Int'l. Conf. Softw. Eng.*, May 2002.
- [26] G. Rothermel, S. Elbaum, A. G. Malishevsky, P. Kallakuri, and X. Qiu. On test suite composition and cost-effective regression testing. Technical Report 03-60-04, Oregon State University, Nov. 2003.
- [27] G. Rothermel and M. J. Harrold. A safe, efficient regression test selection technique. *ACM Trans. Softw. Eng. Meth.*, 6(2):173–210, Apr. 1997.
- [28] G. Rothermel, R. Untch, C. Chu, and M. J. Harrold. Test case prioritization: an empirical study. In *Proc. Int'l. Conf. Softw. Maint.*, pages 179–188, Aug. 1999.
- [29] G. Rothermel, R. Untch, C. Chu, and M. J. Harrold. Test case prioritization. *IEEE Trans. Softw. Eng.*, 27(10), Oct. 2001.
- [30] D. Saff and M. D. Ernst. Reducing wasted development time via continuous testing. In *Proc. Int'l. Symp. Softw. Rel. Engr.*, pages 281–292, Nov. 2003.
- [31] A. Srivastava and T. J. Effectively prioritizing tests in development environment. In *Proc. Int'l. Symp. Softw. Testing and Anal.*, July 2002.
- [32] D. Wells. Extreme programming: A gentle introduction. <http://www.extremeprogramming.org>, Jan. 2003.
- [33] W. Wong, J. Horgan, S. London, and H. Agrawal. A study of effective regression testing in practice. In *Proc. Int'l. Symp. Softw. Rel. Engr.*, pages 230–238, Nov. 1997.