

An Evaluation of Change-Based Coverage Criteria

Marc Fisher II
University of Memphis
Memphis, TN, USA
marc.fisher@memphis.edu

Jan Wloka
IBM Rational Research Lab
Zurich, Switzerland
jan_wloka@ch.ibm.com

Frank Tip
IBM Research
Yorktown Heights, NY, USA
ftip@us.ibm.com

Barbara G. Ryder
Virginia Tech
Blacksburg, VA, USA
ryder@cs.vt.edu

Alexander Luchansky
Vanguard Group, Inc.
Malvern, PA, USA
alexander_luchansky@vanguard.com

ABSTRACT

Various coverage criteria are commonly used to assess the quality of test suites, but achieving full coverage according to these criteria is often impossible or impractical. Our research starts from the popular assumption that a disproportionate number of faults is likely to reside in recently changed code. Based on this assumption, we propose several *change-based coverage criteria* that reflect to what extent changes with respect to a previous program version are exercised by a test suite. In a set of experiments on programs from the SIR repository, we found change-based criteria to reveal faults better than traditional criteria, and to enable the construction of much smaller test suites with similar fault detection effectiveness. We also report on a case study that shows that achieving (near) 100% coverage according to a change-based criterion is both feasible and useful.

1. INTRODUCTION

Software evolves throughout its lifetime as developers adapt applications to changing requirements, add new functionality, and fix bugs. To ensure that their changes do not have unintended consequences, developers typically rely on a suite of unit or regression tests. The hope is that, by executing the test suite after each edit, faults are exposed early when understanding and fixing the problem is relatively easy as the changes are still fresh in the developer's mind. The effectiveness of this approach hinges on the quality of the test suite, (i.e., its ability to expose faults). Following Miller and Maloney [14], various notions of code coverage have been proposed as a proxy for test suite quality, including statement coverage, branch coverage, method coverage, and others [15, 19, 34]. The basic assumption is that a test suite is likely to be effective at revealing faults if it exercises the code where the fault is located. Therefore, increased code coverage is expected to correlate with more revealed faults.

Unfortunately, achieving 100% coverage according to a traditional code coverage criterion is often neither possible nor practical. Achieving full coverage becomes impossible when applica-

tions contain unreachable code, which may arise due to programmer carelessness, or when the implementation of a feature is not yet complete. Furthermore, only a limited amount of time may be available for writing tests, and achieving full coverage may not be practical because it is often quite difficult to construct a test that exercises a particular area of the code. Therefore, organizations often select some arbitrary coverage level (e.g., 70%) as their target [4]. In our opinion, imposing an arbitrary coverage level target is problematic because there is no guarantee that programmers will focus their efforts on writing tests that cover the most fault-prone parts of an application.

It has long been understood that changed code requires more attention than code that has been stable for a long time and a variety of program analysis techniques have been developed to analyze the effect of changes. For example, regression test selection techniques identify those tests whose behavior may have been changed as a result of code changes [23], change impact analysis techniques can identify the parts of a program that are potentially affected by a set of changes [6], change classification techniques identify the subset of a programmer's changes that is responsible for the failure of a given test [28], safe-commit analysis determines changes that can be committed safely to a version control repository [32], and test suite augmentation techniques provide guidance for creating new test cases that specifically target a program's changed behavior [2, 21, 26]. All of these techniques involve some form of program analysis, starting from the changed code fragments, with the objective of helping the programmer improve or preserve code quality.

In this paper, we suggest that a test suite's quality should be measured directly by its ability to cover changes. To this end, we propose a number of *change-based coverage criteria* that reflect to what extent the *changes* with respect to a previous version of the program are exercised by a test suite. In particular, we define novel *changed statement*, *changed method*, *changed branch*, and *atomic change* coverage criteria. The premise of our research is that the change-based coverage achieved by a test suite should correlate well with its fault-finding ability because faults are likely to be concentrated in recently modified code. Furthermore, as the number of changed statements (or methods or branches) is typically much smaller than the total number of statements (or methods or branches), change-based coverage criteria are likely to be more *practical* than traditional coverage criteria.

We conducted a set of experiments to compare the fault-finding ability of test suites constructed according to the 4 change-based coverage criteria mentioned above to that of test suites constructed according to 3 traditional (*statement*, *method*, and *branch*) coverage criteria. To evaluate the relative effectiveness of these 7 criteria, we considered 12 version pairs of 3 existing Java applications from the

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to publish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.
Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$10.00.

SIR repository [7], in which faults were seeded. For each version pair and each criterion, we then determined the coverage achieved and number of faults exposed for many subsets of the originally provided test suite. From this information, we computed a correlation between coverage and exposed faults using the Kendall τ method [10]. Our results indicate that change-based coverage criteria correlate more strongly with exposed faults than non-change-based coverage criteria on all but one of the version pairs.¹ We did not find any change-based coverage criterion to be clearly superior to all others. However, for 10 of the 12 version pairs, one of the *changed branch* and *changed statement* criteria correlates most strongly with exposed faults, and in the remaining two cases, one of these criteria is very close to the winner. In our opinion, these results experimentally confirm that it is a good idea to focus testing effort on changed code fragments, and that it is particularly important to cover changed branches and statements.

In a second set of experiments, we evaluate the tradeoff between test suite size and fault-detection effectiveness, for test suites created according to the different coverage criteria. To this end, we minimized the existing test suite for each version pair under test in a way that preserved coverage according to each criterion and measured the fault-detection effectiveness of the resulting minimized test suites. Our results indicate that minimizing a test suite according to the change-based criteria can decrease suite size significantly, without decreasing the number of faults exposed by much. For example, we found that test suites minimized according to the *changed statement coverage* criterion are on average between 36% and 89% smaller than the test suites minimized according to the *statement coverage* criterion. However, these smaller test suites were equally effective at revealing faults on 9 of the 12 version pairs under consideration, and nearly as effective on 2 additional version pairs. These results again suggest that, when limited time is available for writing additional tests, priority should be given to covering changed code fragments.

We also conducted a case study in which we investigated, for 7 version pairs of two of the subject programs, how practical it is for programmers to write additional tests that achieve 100% coverage according to the *atomic change* criterion. We found this task to be feasible in terms of effort, and surprisingly useful for pointing programmers at untested changes that might otherwise be overlooked. In particular, the additional tests written revealed 3 previously unexposed seeded faults and 2 real faults not previously reported (an arithmetic error and a situation where one of the programs throws a `NullPointerException` when given a malformed input), thus providing some evidence that the use of change-based test adequacy criteria may help prevent errors and improve code quality.

In summary, the contributions of this paper are as follows:

- We define several new change-based code coverage criteria.
- We compare the effectiveness of these criteria to that of several traditional code coverage criteria. Empirical results indicate that change-based criteria reveal faults better than traditional criteria, and that they can be used to construct much smaller test suites with similar fault detection effectiveness.
- We report on a case study that shows that achieving (near) 100% coverage according to a change-based coverage criterion is both feasible and useful.

These results provide some quantitative evidence to confirm the popular belief that it is worthwhile to focus testing effort on changed code fragments, and suggest that change-based coverage criteria can serve as a practical measure of test suite adequacy.

¹ In the one remaining case, a non-change-based criterion was only slightly better.

```

class Bar {
void zap() {
    System.out.println("Bar.zap()");
}
}
class Foo extends Bar {
void zip(String s) {
    System.out.print(s.length() + ": ");
    System.out.print("Foo.zip()");
    if (s.length() == 1 || s.length() == 2 ) {
        System.out.println("1");
    } else {
        System.out.print("--");
        System.out.println("2");
    }
}
}
void zap() {
    System.out.println("Foo.zap()");
}
}

```

Figure 1: Original and Edited Version of the Example Program with Edits and Atomic Changes (Textual edits are shown with shaded boxes. Atomic changes (CM, LC) are indicated by framed boxes. The original version of the program excludes all shaded code fragments, and edited version can be constructed by adding these fragments to the program.)

The remainder of this paper is organized as follows. Section 2 defines the change-based coverage criteria. Our empirical evaluation is presented in Section 3. A case study involving writing new tests to cover the changes made during edits of *NanoXML* and *JTopas* is presented in Section 4. Threats to validity are discussed in Section 5. Finally, related work is discussed in Section 6 and conclusions are presented in Section 7.

2. COVERAGE CRITERIA

Code coverage is a white-box testing technique to measure how much of a program is exercised by a test suite. Each test in the suite uses assertions to specify expected program behavior. A run of a test suite tries to verify that the program works as expected. Each run computes the degree to which the program is covered by tests as measured by a specific coverage criterion, which is defined in terms of a measurable unit, such as a method, statement, or branch. In this section, we give an overview of all the coverage criteria compared in our empirical study.

2.1 Illustrating Example

Figure 1 shows two versions of a small example program that we will use to illustrate several traditional and change-based coverage criteria. The original version contains the lines of code that lie outside of the shaded areas. The edited version can be obtained by adding all code in the shaded boxes. *Edit 1* adds a new print statement, *Edit 2* adds a test to the `or` condition in the `if` statement, *Edit 3* adds a new print statement to the `else` branch, and *Edit 4* adds a method `Bar.zap()` that overrides `Foo.zap()`. The four atomic changes computed from the textual edits and labeled on the right side of Figure 1 are defined in Section 2.3. For each coverage criterion, we discuss how a test suite can be constructed that attains full coverage of the program, summarizing each test suite simply by the set of method calls that it performs.

2.2 Traditional Coverage Criteria

These criteria are measured with respect to the complete program. Method, statement and branch coverage are commonly used.

Method Coverage.

Method (or function) coverage [34] is commonly used to measure which methods (functions) in the program are invoked during test execution. It indicates whether at least one statement within a method is executed at least once during a test suite run. Method coverage does not differentiate the individual statements actually executed. If we consider the edited version of the example program, 100% method coverage can be achieved by running a test suite containing the following method calls:

```
new Foo().zip(""); new Foo().zap();
new Bar().zap();
```

Statement Coverage.

Another common coverage measure is statement (or line) coverage [15]. Statement coverage reports all statements that are invoked at least once during a test suite run. Statement coverage subsumes method coverage, (i.e., a test suite with 100% statement coverage will have 100% method coverage). A test suite with the following method calls reports 100% statement coverage in the example program

```
new Foo().zip(""); new Foo().zip("a");
new Foo().zap(); new Bar().zap();
```

The extra call to `new Foo().zip("a")` exercises the statements on the then branch of the `if` statement, which were not covered by the previous test suite for method coverage.

Branch Coverage.

Branch coverage [15] requires that each outcome (e.g., *true* and *false*) of each branch in a logical expression controlling an `if`, `loop`, or `switch` statement be covered. For practical reasons (we use instrumentation tools that operate at the bytecode level), we assume that composite logical expressions are decomposed into a sequence of lower-level branches. For example, using short circuit execution, full branch coverage of a logical `or` expression `a or b` can be obtained with: (i) `a ==T`, (ii) `a==F, b==T` and (iii) `a==F, b==F`. Note that the last case is required to be covered, but the settings `a==T, b==F` and `a==T, b==T` are not.

As usually defined, branch coverage subsumes statement coverage [19]. To maintain this constraint, in our implementation we require branch coverage to cover all methods in addition to all decisions within these methods. A set of method calls that achieves 100% branch coverage on the example is:

```
new Foo().zip(""); new Foo().zip("a");
new Foo().zip("ab"); new Bar().zap();
new Foo().zap();
```

2.3 Change-based Coverage Criteria

The basis for the change-based coverage criteria used in our empirical study is a previously developed change impact analysis implemented in the tool CHIANTI [25, 20, 31]. This analysis computes an abstract representation of a textual difference between two program versions and decomposes it into a set of atomic (or smallest) changes to a program. The resulting representation enables a classification of different kinds of changes and their dependences, making textual edits amenable to program analysis. Various change categories are supported, such as change method (CM), add field (AF), and lookup change (LC) (i.e., a change to dynamic dispatch) [20]. However, in this paper we use a simplified version of the

change model including only change method (CM) changes and lookup (LC) changes.

Atomic Change Coverage.

The *atomic change coverage* criterion reflects the CM and LC atomic changes that are exercised by a test suite. CM changes correspond to methods whose code has changed. CM changes also summarize other changes; for example, adding a new instance field to a class implicitly changes the constructor method for that class to initialize an additional field. Note that we generate one CM change per method regardless of the number of changed statements within the respective method's body, (e.g., CM1 in the example represents *Edits 1 – 3* as a single CM change).

A lookup change (LC) represents the effect of an edit on dynamic dispatch. It is represented as a pair of the form $\langle C, f() \rangle$ indicating that the behavior of invoking a method $f()$ on an object of runtime type C has changed. Many kinds of edits may alter the existing dynamic dispatch behavior of a Java program, such as adding the overriding method `zap()` to class `Foo` in our example (LC1, LC2), or changing method visibility from *private* to *public* [20, 25, 31]. A test suite with the following method calls covers all CM and LC changes and achieves 100% atomic change coverage on the example program:

```
new Foo().zip(""); new Foo().zap();
new Bar().zap();
```

Changed Method Coverage.

The *changed method coverage* criterion reports how many of the CM changes are covered by a test suite. It uses the same approach as atomic change coverage but is restricted to CM changes. As such, changed method coverage is subsumed by atomic change coverage. The following method calls within a test suite result in 100% changed method coverage:

```
new Foo().zip(""); new Foo().zap();
```

Changed Statement Coverage.

A test suite's *changed statement coverage* reflects the set of statements that were changed or added since the previous version and exercised by at least one test in the suite. A test suite with the following method calls reports 100% changed statement coverage:

```
new Foo().zip(""); new Foo().zap();
```

A change impact analysis at the level of individual statements (or lines) implemented by JDIFF [1] is used to compute changed statement coverage. Two different versions of a program are compared and both differences and correspondences at statement-level are identified to capture all the changed statements in a Java program.

Changed Branch Coverage.

There are multiple possible ways to define a changed branch coverage criterion. For example, one could focus only on those branches where the condition has changed, or one could also include branches where the body has changed. Since we are not aware of any existing definition or approach for identifying a set of changed branches within an application, we define a simple, conservative notion of changed branch coverage based on the set of CMs identified by CHIANTI. Our *changed branch coverage* criterion requires that all branches within each changed method are covered, even if none of the edits affected the corresponding control statement. In addition, the criterion requires that all CM changes are covered, in order to ensure that changed branch coverage subsumes changed statement coverage. According to this criterion, a

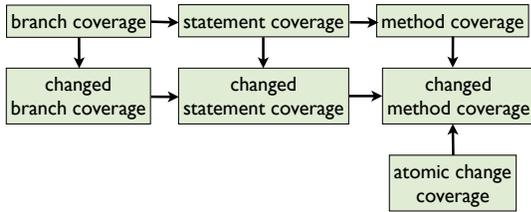


Figure 2: Subsumption relationships between the criteria.

test suite with the following method calls achieves 100% changed branch coverage:

```

new Foo().zip(""); new Foo().zip("a");
new Foo().zip("ab"); new Foo().zap();
  
```

2.4 Subsumption Relationships

Figure 2 shows the subsumption relationships that hold between the coverage criteria discussed in this section. The presence of an arrow from a criterion A to a criterion B in the figure means that 100% coverage according to criterion A implies 100% coverage according to criterion B . Note that atomic change coverage is not subsumed by any other coverage type as it is the only criterion to consider the coverage of (changed) method dispatch behaviors.

3. EXPERIMENTAL EVALUATION

In order to better understand the properties of the different coverage criteria, we conducted a study that compares these criteria by answering the following research questions:

- RQ1** How well do the different coverage criteria predict the relative quality of test suites?
- RQ2** How do test suites constructed according to the different coverage criteria compare in terms of size and number of regression faults exposed?

3.1 Coverage Criteria

Our empirical study compares the effectiveness of the seven different code coverage criteria defined in Section 2: *method* coverage, *statement* coverage, *branch* coverage, *changed method* coverage, *changed statement* coverage, *changed branch* coverage and *atomic change* coverage. The first three of these are traditional code coverage metrics in common use. To identify the set of changed statements empirically for each program version, we used the JDIFF tool [1]. COBERTURA, version 1.9 [8] was used to run each test suite and to measure method, statement, branch, changed statement, and changed branch coverage. We used JUNITMX [31] to identify atomic changes, run each test suite, and calculate changed method and atomic change coverage as well as the faults exposed by each suite.

As defined in Section 2, changed statement coverage subsumes changed method coverage and changed branch coverage subsumes changed statement coverage. However, JDIFF computes changes by analyzing bytecode, while CHIANTI computes changes by analyzing the source code. It is possible for unchanged source code within a method to produce different bytecode when compiled (e.g., when a call site is resolved to a different method definition) or for minor changes in source code to result in identical bytecode. This can lead to differences in the set of methods with changes identified by JDIFF and CHIANTI, breaking the subsumption relationships discussed in Section 2.4. We experimentally determined that these

benchmark	version	tests		faults	
		all	used	all	rev
<i>JMeter</i>	v0/v1	77	77	20	5
	v1/v2	80	80	21	5
	v2/v3	78	77	20	8
	v3/v4	78	75	13	1
	v4/v5	97	80	13	2
<i>JTopas</i>	v0/v1	126	126	10	7
	v1/v2	128	128	12	6
	v2/v3	209	206	17	8
<i>NanoXML</i>	v0/v1	214	214	7	7
	v1/v2	214	214	7	6
	v2/v3	216	216	11	10
	v4/v5	216	216	9	9

Table 1: Benchmark Details

minor discrepancies did not significantly affect the results reported in this section (see Section 5).

3.2 Benchmarks

To evaluate our technique, we used three existing Java applications, *JMeter*, *JTopas*, and *NanoXML*, available from the SIR repository [7]. The SIR repository contains between four and six versions for each of these applications, with an associated test suite for each version, and with associated seeded faults for most versions. Since this study is focused on changes to the software, we used successive version pairs of each application from the SIR repository as shown in column 2 of Table 1, resulting in 12 different version pairs.²

Column 3 of Table 1 indicates the number of tests available for the most recent version of each version pair. COBERTURA required that each test be run in isolation from the rest of the tests in the test suite; however, some tests for *JMeter* v3, v4, and v5 could only be successfully executed when run as part of the entire test suite. Therefore, we removed these tests from the test suites that we used for this study; column 4 shows the number of tests actually used. The faults columns of Table 1 show the total number of seeded faults for each version pair (all) and the number revealed by the pool of tests we used for each version (rev). Table 2 provides the number of coverable elements (all) and the number of those elements covered by the test suite used (cov) for each of the test coverage criteria we considered.

3.3 Experimental Methodology

For each version pair, V_{old} and V_{new} , we created faulty versions of V_{new} by enabling a single fault within V_{new} . Then, we computed coverage and fault detection matrices by executing each test case t in the test suite on each faulty version of V_{new} as well as on the correct version of V_{new} . For each t , we created coverage and fault detection matrices that indicated the coverable elements covered by t and the faults exposed by t , respectively. We considered an element covered by the execution of t if for at least one version of V_{new} (correct or faulty) that element was covered. A fault f was exposed by test t if t failed on the version of V_{new} with that fault f enabled. The coverage values for a test suite were calculated by unioning the matrix entries for their constituent test cases.

²Version 4 of *NanoXML* did not include any seeded faults so we were unable to use the v3/v4 version pair of this application

benchmark	version	meths		stmts		brchs		chg meths		chg stmts		chg brchs		atomic chgs	
		all	cov	all	cov	all	cov	all	cov	all	cov	all	cov	all	cov
<i>JMeter</i>	v0/v1	3082	1222	16993	6865	7042	2059	897	337	5761	2451	2701	641	1974	534
	v1/v2	3033	1192	16948	7036	7186	2029	281	119	2408	1276	1156	236	1209	401
	v2/v3	3650	1591	18412	7720	8541	2598	1530	874	8605	4760	4121	1320	5294	2483
	v3/v4	3767	1666	19038	8050	8635	2679	370	145	3163	1462	1062	253	1046	331
	v4/v5	3849	1701	20519	8571	8677	2735	159	74	2069	766	826	204	247	101
<i>JTopas</i>	v0/v1	199	164	796	683	625	434	30	27	185	157	60	46	55	49
	v1/v2	213	178	878	744	726	490	46	41	223	175	154	90	79	70
	v2/v3	490	349	2181	1620	1692	1025	360	245	1404	984	1242	732	797	549
<i>NanoXML</i>	v0/v1	120	98	925	717	552	376	68	52	645	471	416	266	110	85
	v1/v2	177	106	1100	745	625	394	126	69	684	428	411	260	242	97
	v2/v3	214	126	1391	921	810	487	78	54	651	446	334	204	187	115
	v4/v5	229	134	1509	981	876	508	64	50	590	375	420	279	84	59

Table 2: Criteria Details

3.4 RQ1

The goal of RQ1 was to determine which of our coverage criteria most accurately measures the relative quality of different test suites with respect to exposing regression faults. To answer this question, we needed to study the effect of varying coverage on fault exposure by measuring the coverage and fault exposure on several test suites for each VUT (i.e., version pair under test). Since we only were supplied with a single test suite for each version, we needed to construct additional test suites. This was accomplished by randomly selecting 100,000 different subsets of the original test suite for each VUT. Note that we only needed to run each test in the original test suite once. After randomly constructing new test suites, we used the previously collected coverage and fault exposure data to calculate the coverage and exposed faults for each test suite.

Table 3 shows the correlation coefficients computed for coverage vs. fault detection for each of the coverage criteria on each of the VUTs. The coefficients were computed using the Kendall τ method [10]. Kendall τ is a standard test that computes the correlation between the ranks of the values for two variables. Kendall τ is 1 if the ranks of the values for the two variables exactly match, -1 if the ranks of the values are in exactly the opposite order, and near 0 if the variables are not related. The correlation coefficients of different coverage criteria can be compared to determine which criterion more accurately predicts the relative strengths of the different test suites at exposing faults. Specifically, if criterion C_1 correlates more strongly with fault detection (i.e., has a higher correlation coefficient) than criterion C_2 , then C_1 is the better criterion for predicting which of two test suites is more likely to reveal regression faults. It should be noted, however, that this says nothing about whether a particular test suite with high C_1 coverage is more likely to reveal regression faults than another test suite with equally high C_2 coverage.

For 7 of the 12 version pairs, changed branch coverage showed the strongest correlation with fault exposure. These seven cases included all version pairs of *JMeter* and two of the three *JTopas* pairs. Looking at the cases where changed branch coverage was not the winner, the primary difference seems to be the percentage of all of the branches that were identified as changed. For example, on *JTopas* v2/v3, where changed statement coverage won, 73% of all of the branches were identified as changed. In contrast, on *JTopas* v1/v2, where changed branch coverage won, only 21% of the branches were identified as changed. In general, the five runs where changed branch coverage was not the best had a higher proportion of changed branches (and changed methods and changed statements) than the seven runs where changed branch coverage

was the best (although there was some overlap in the ranges). This suggests that the size of the change as a proportion of the entire application is an important component in determining the effectiveness of using a change-based coverage criterion. This is not unexpected as the main purpose of change-based criteria is to focus attention on testing the changed portions of the application, and, in cases where large portions of the code have changed, it is reasonable to expect that such a focus would be less effective. However, it should be noted that in four of the five cases where changed branch coverage was not the best, another change-based criterion performed the best, suggesting that even with large sets of changes a change-based criterion could be useful in guiding testing effort.

Another way to look at these results is by comparing the correlation coefficient for traditional coverage criteria to those for their change-based counterparts. Comparing the results for *method* coverage to those for *changed method* coverage, we see higher correlation coefficients for the latter on 7 of the 12 version pairs. The *changed statement* coverage criterion outperformed the traditional *statement* coverage criterion in 10 of 12 cases, and the *changed branch* criterion did better than the traditional *branch* coverage criterion in 9 of 12 cases.

These results suggest that by focusing attention on changes, we can better indicate the relative quality of test suites in exposing regression faults. This result may not seem surprising as the faults are in the changed portions of the applications; therefore, executing these portions of the programs is likely to execute the faulty statements. However, it is quite possible that the context in which a fault is executed is just as important as executing the fault itself. Therefore, a criterion that emphasizes exercising more of the program may be a better indicator of the regression fault detection effectiveness of a test suite.

Summary.

We found that:

- change-based criteria provide better predictions about the relative regression fault detection exposure of test suites than non-change-based criteria;
- of the change-based criteria, changed branch coverage generally provided the best prediction.

3.5 RQ2

Our second research question focused on evaluating the trade-off between test suite size and fault-detection effectiveness, for test suites created using the different criteria. The ideal evalua-

benchmark	version	meths	stmts	brchs	chg meths	chg stmts	chg brchs	atomic chgs
<i>JMeter</i>	v0/v1	0.5180	0.5237	0.5241	0.5075	0.5203	0.5439	0.5073
	v1/v2	0.4762	0.4685	0.4798	0.5438	0.5284	0.5922	0.4801
	v2/v3	0.5350	0.5326	0.5415	0.5392	0.5336	0.5449	0.5211
	v3/v4	0.4386	0.4417	0.4451	0.5073	0.4389	0.5416	0.4359
	v4/v5	0.5183	0.5200	0.5276	0.5740	0.5669	0.6767	0.5884
<i>JTopas</i>	v0/v1	0.6767	0.7081	0.6631	0.7443	0.7343	0.7826	0.7376
	v1/v2	0.6110	0.6783	0.6931	0.6637	0.7516	0.7968	0.6903
	v2/v3	0.6487	0.6728	0.6634	0.6555	0.6799	0.6453	0.6472
<i>NanoXML</i>	v0/v1	0.5450	0.6516	0.5917	0.5078	0.7422	0.6534	0.5078
	v1/v2	0.7509	0.6792	0.7040	0.7353	0.6821	0.7095	0.7529
	v2/v3	0.8093	0.7693	0.7917	0.7313	0.7792	0.7736	0.7976
	v4/v5	0.5505	0.5697	0.5430	0.3286	0.7254	0.5137	0.3286

Table 3: Correlation Coefficients (bold indicates best value for row, all values were statistically significant ($p < 0.05$))

object	version	meths		stmts		brchs		chg meths		chg stmts		chg brchs		atomic chgs	
		size	rev	size	rev	size	rev	size	rev	size	rev	size	rev	size	rev
<i>JMeter</i>	v0/v1	21	3.35	36	3	36	3	15	3.1	27	3	22	3.33	17	3.2
	v1/v2	23	5	38	5	36	5	11	5	23	5	13	5	16	5
	v2/v3	28	8	37.8	8	34	8	27	8	33.8	8	28	8	30	8
	v3/v4	26	1	36	1	32	1	9	1	22	1	10	1	9	1
	v4/v5	26	2	36	2	33	2	11	1.05	20	2	14.7	1.65	12	2
<i>JTopas</i>	v0/v1	15	7	21	7	21	7	6	5.2	14.1	7	6	7	6	5.3
	v1/v2	15	6	22	6	22	6	8	6	12	6	6	5	10	6
	v2/v3	30	4.5	45	8	46.65	8	20	4.1	35	7.5	19.6	7	20	4.15
<i>NanoXML</i>	v0/v1	15	5.7	25.5	7	21	6.45	5	3.55	11.95	6.25	10	5.55	5	3.8
	v1/v2	16	4.7	25.5	5.25	21	4.3	10.2	4.65	14.55	5.25	11	4.25	10	4.7
	v2/v3	17	6.5	29	10	24	9	10	3.6	22	10	15	9	13	6.5
	v4/v5	17	7.4	29	9	23	7.7	3	4.25	10.55	6.35	7	7	3	3.5

Table 4: Average Size of and Faults Revealed by Minimized Test Suites

tion mechanism would have been to generate a variety of adequate test suites for each of the criteria and compare them. Unfortunately, doing so in a fashion that allows fair comparisons of the criteria is expensive and difficult. Therefore, we approached this problem by minimizing the existing test suite for each VUT using each of the coverage criteria.

Since the goal of our minimization was to select subsets of the original test suite without sacrificing coverage, an optimal or near-optimal minimization technique was not required. Instead, we implemented a simple greedy minimization technique that allowed us to readily generate different test suites with no completely redundant tests. Our minimization technique begins with an empty test suite, S . Then, we incrementally add a test to S by randomly selecting a test for the VUT that increases the coverage of S . This step is repeated until the coverage achieved by S equals the coverage of the entire test suite for the VUT. We then iterate through the tests in S , identifying and removing any test whose removal does not decrease the coverage of S . While this approach does not guarantee a minimal test suite, in practice we found that it effectively generated small test suites that varied little in size and included different tests.

For our experiments, we generated 20 unique test suites for each criterion, VUT pair. For each of these test suites, we calculated the number of faults exposed and the test suite size. The averages of these two metrics across the 20 test suites is presented in Table 4.

Although method coverage minimized suites revealed more faults than changed method coverage minimized suites on eight of the version pairs, this difference was small, only exceeding one fault

on four of the version pairs. Interestingly, this increase in fault exposure comes at considerable expense. The test suites for method coverage tended to be significantly larger than the test suites for changed method coverage, ranging from 1.04 times as large to 5.67 times as large, and exceeding 2.0 times as large on 7 of the 12 version pairs. The same pattern was found when comparing statement coverage minimized suites with changed statement coverage minimized suites and branch coverage minimized suites with changed branch minimized suites, although the changes in number of exposed faults and in test suite size tended to be smaller. These findings show that a significant benefit in efficiency can be gained by focusing on testing changes while missing few regression faults.

Comparing the different change-based criteria to each other, we see that changed statement coverage minimized suites reveal more faults and contain more tests on average than the other change-based criteria minimized suites and changed method coverage minimized suites generally reveal fewer faults and contain fewer tests than the other minimized suites. There is not a consistent difference in ranking between changed branch minimized suites and atomic change minimized suites. Thus, for change-based coverage criteria, there is the usual trade-off between fault-exposure potential and test suite size.

Summary.

We found that:

- test suites constructed according to change-based criteria found approximately the same number of faults on average, but

were significantly smaller than test suites constructed according to corresponding non-change-based criteria;

- there is a trade-off between fault-exposure potential and test suite size for test suites constructed according to different change-based criteria.

4. CASE STUDY

In this section we report on a case study investigating the feasibility of achieving 100% change coverage using the atomic change criterion. We chose this criterion for three reasons:

- It provides information (computed using program analysis) about semantic changes in object-oriented programs that can guide the developer in writing new tests (e.g., the addition of a method resulting in a new target of a dynamic dispatch);
- Since this criterion tracks method-level changes, there are fewer changes to cover than, for example, with statement coverage, making 100% coverage a more feasible goal; and
- We have a long-standing interest in analyzing the impact of changes in a object-oriented setting [20, 25, 31].

To demonstrate the effectiveness of the atomic change criterion, we ran the test suites for each changed version of two SIR benchmarks –*NanoXML* and *JTopas*– and then tried to write additional tests to cover all changes reported as uncovered. Although our results are anecdotal, we believe they demonstrate that the atomic change criterion aided the construction of new tests and thereby, the achievement of near 100% change coverage.

Quantitative Results.

Table 5 summarizes our quantitative results. For each version pair the table shows: (i) the number of uncovered changes, (ii) the number and LOC of additional tests written, (iii) the number and LOC of the new input files used by the additional tests,³ and (iv) the number of additional seeded and real (i.e., non-seeded) faults exposed by the additional tests.

For example, for the v2/v3 version pair of *NanoXML*, 13 additional tests were written, comprising a total of 135 lines of source code, in order to cover CM and LC changes reported as uncovered by our tool. These additional tests required 3 new input files, which comprised a total of 21 lines of XML. No additional seeded faults were exposed by the tests, but they did uncover a real fault that was accidentally introduced by the developers as they made their changes!

For our benchmarks, newly written tests exercising the uncovered CM changes usually also covered any reported LC changes. Hence, the use of the changed method coverage criterion would have produced similar results in terms of the required effort and the number of exposed faults. To a large extent this is because neither of the benchmarks under study was written in a very object-oriented style, therefore, we suspect that different subject programs with more complex class hierarchies might have yielded different results.

Test Creation Process.

Initially, we were unfamiliar with the code of these benchmarks. After some initial experience with trying to cover methods that were previously not exercised by the test suite, we quickly realized

³ Only *NanoXML* required the creation of additional input files, containing XML code.

that developing new test cases “from scratch” can be quite challenging. As a result, we converged on a process for deriving new tests from existing ones. Specifically, we search for a covered method m containing a call to a changed method m' that we need to cover (if we are unable to do this, we try to find a covered method m that indirectly calls m'). Presumably, the call to m' is not executed because a control condition is not satisfied (e.g., the call to m' is on an unexecuted else-branch of an if-statement). In such cases, the challenge is to understand how—by changing the program input—we can persuade the program to select the branch containing the call to m' . In the absence of automated tool support for this task, we had to inspect the program state in the debugger and understand how the conditions that guard the call to m' depend on program inputs or constant values. Following this approach, we were able to cover all but a few of the atomic changes with relatively little effort.

The use of a change-based criterion helped us to focus on new tests covering the edited portions of the program. In contrast, to achieve 100% coverage with a non-change-based criterion, we would have introduced unnecessary tests for portions of the program unaffected by the edit. Additionally, by reporting changes at the method-level, the atomic change criterion suggested a natural structure for the additional tests that statement-based criteria did not.

We estimate that the amount of effort involved in writing the additional tests was in the order of 2-8 hours per version pair; however, we conjecture that the original developers could have performed this task in a fraction of the time, given their familiarity with the code.

JTopas.

Achieving atomic change coverage of the versions of *JTopas* resulted in the exposure of three seeded faults not covered by the original test suite and discovery of a non-seeded fault in one version!

For *JTopas*, most of the changes were of a trivial nature and could be covered with relatively little effort. For example, in the *JTopas* v1/v2 version pair, an uncovered getter method `PluginTokenizer.getKeywordHandler()` was tested by adding a simple assertion to a test that already invoked the corresponding `setKeywordHandler()` setter method.

A more interesting case in *JTopas* v1/v2 was a new uncovered method `AbstractTokenizer.getTextUnchecked(int, int)`, which is similar to an existing method `AbstractTokenizer.getText(int, int)`. The new `getTextUnchecked()` method does not check to make sure that its parameters are within range, but if they are not, the JavaDocs report that “*a java.util.IndexOutOfBoundsException may occur, or uninitialized data may be retrieved*”. Covering this method is a bit involved, because it is not intended to be invoked directly from a test. However, the `getText()` method is invoked by a method `current()`, a method in the same class, which is invoked by several tests, including `TestInputStreamTokenizer.testLineCounting()`. We covered `getTextUnchecked()` by adding: (i) a `currentUnchecked()` method, which is similar to `current()` but calls `getTextUnchecked()` instead of `getText()`, and (ii) a test `testLineCounting2()` that is identical to `testLineCounting()` except that it calls `currentUnchecked()` instead of `current()`. Interestingly, this added test exposed one of the seeded faults for *JTopas* that was not exposed by the original test suite, thus illustrating how the use of a change-based coverage criterion can help improve software quality.

benchmark	version	uncovered changes	added tests		added inputs		#faults exposed	
			#	LOC	#	LOC	seeded	non-seeded
<i>JTopas</i>	v0/v1	5	3	91	0	0	1	0
	v1/v2	7	2	74	0	0	0	0
	v2/v3	106	31	104	0	0	2	1
<i>NanoXML</i>	v0/v1	18	9	114	6	33	0	0
	v1/v2	86	10	420	6	63	0	0
	v2/v3	28	13	135	3	21	0	1
	v4/v5	16	9	127	9	103	0	0

Table 5: Quantitative Results of Case Study

Another interesting case occurred in the *JTopas* v2/v3 version pair. While adding new tests that exercise the uncovered changes, we discovered a bug in the newly added `Token.getEndPosition()` method (i.e., a non-seeded fault). This method is currently implemented as:

```
public int getEndPosition() {
    return getLength() - getStartPosition();
}
```

but it should be:

```
public int getEndPosition() {
    return getLength() + getStartPosition();
}
```

This newly added method was not covered by any of the existing test cases, and the use of atomic change coverage as a test adequacy criterion led us directly to this non-seeded fault.

NanoXML.

Achieving atomic change coverage of the *NanoXML* benchmark versions resulted in the exposure of a non-seeded fault. Moreover, many of the uncovered changes revealed the inadequacy of the test suite in testing newly introduced code for handling exceptions, as well as some instances of poor Java coding practice using `finalize()`. Although 100% coverage was not possible due to the latter cases, the quality of the original test suite was strengthened by providing tests of exceptional conditions, which are often ignored.

NanoXML is an XML parser. Many of the tests in its existing test suite apply the parser to an XML input file and execute some assertions that compare the actual output against the expected output. Several of the tests apply the parser to malformed XML and check whether or not an expected exception is thrown. Most of the uncovered changes for *NanoXML* correspond to error handling for various forms of malformed XML that the original test suite did not exercise. Writing additional tests to cover such functionality involves creating a new malformed XML input file and creating a test case with the appropriate assertions and/or exception handlers. This was a bit challenging because we were not familiar with either the *NanoXML* source code or with the finer details of XML syntax.

As a specific example of the kind of new tests needed, consider the *NanoXML* v0/v1 version pair, which includes a method `XMLUtil.errorInvalidInput(IXMLReader, String)` containing some error handling code not exercised by the existing test suite. After some experimentation, we determined that this method will be executed if a line

```
<!ENTITY % extParamEntity "blah">
```

in the input file associated with an existing test `testParser1_vw_v1_17` is changed to the following:

```
<!ENTITY >> extParamEntity "blah">
```

Executing *NanoXML* on the resulting modified input file resulted in an `XMLException` being thrown with an associated message

```
XML Not Well-Formed at Line 3:
Invalid input: e
```

We added a new test `testParser1_vw_v1_17B` that parsed the resulting modified input file, caught the `XMLException` and contained an assertion to check the error message. In our opinion, this is precisely the kind of test the programmer should have written to cover this functionality.

In the course of writing a test to exercise the uncovered method `XMLElement.getFirstChildNamed(String, String)` in the *NanoXML* v2/v3 version pair, we encountered a non-seeded fault (i.e., a bug) that causes a `NullPointerException` to be thrown if the parser is invoked on certain malformed XML input files. After some investigation, we tracked the fault down to an assignment

```
found &= (name == null);
```

in `XMLElement.getFirstChildNamed(String, String)` that should have been

```
found &= (namespace == null);
```

This provides further evidence that the use of atomic change coverage as a test adequacy criterion can help developers identify faults that would be missed otherwise.

Our final discussion concerns uncovered `finalize()` methods in the *NanoXML* v4/v5 version pair. These finalizers appear in newly added classes `PIReader` and `CDATAREader` in an attempt to ensure that reader-objects wrapped by these classes are properly discarded. This is an erroneous coding practice because the Java runtime provides very weak guarantees about the execution of `finalize()` methods, and the use of finalizers for time-critical and resource-critical operations is strongly discouraged.⁴ In this case, the preferred solution is to replace the finalizers with explicitly invoked clean-up methods. We considered debugging the subject programs to be outside the scope of our study and therefore left these methods uncovered, achieving only near 100% coverage.

Summary.

We found the use of atomic change coverage as a test adequacy criterion feasible in terms of the amount of effort involved, and useful for pointing programmers at untested changes that might otherwise be overlooked. **The case study revealed 3 previously unexposed seeded faults and 2 real faults, thus providing some evidence that using atomic change coverage as a test adequacy criterion may help prevent errors and improve code quality.**

⁴ In his seminal book, Bloch [5, page 27] writes that “Not only does the language specification provide no guarantee that finalizers will get executed promptly; it provides no guarantee that they’ll get executed at all.”. Regarding the use of finalizers to release nonmemory resources, he goes on to say that “...it is a grave error to depend on a finalizer to close open files because open file descriptors are a limited resource. If many files are left open, because the JVM is tardy in executing finalizers, a program may fail because it can no longer open files.”.

5. THREATS AND LIMITATIONS

There are a wide range of threats and limitations relevant to our research. The primary threats to validity are external, affecting the generalizability of our results. These include the choice of benchmarks and the versions, test suites and seeded faults for these benchmarks. Any of these choices may not be representative of what is found within the wider software development community and, therefore, our results may not generalize. To mitigate these threats, we used benchmarks from the SIR repository that have been used in the evaluation of a wide variety of testing methodologies. The particular selected benchmarks were chosen based on the ability of the tools we were using to handle them and to cover a range of different applications.

Our choice of benchmarks and study design also limits the applicability and interpretation of our results. In particular, we used sampled subsets of the tests provided for these benchmarks as proxies for different possible test suites. In practice, it is possible that test designers creating tests for a particular test coverage criteria would create different types of test suites than those studied. Additionally, our goal was to study the usefulness of different coverage criteria in finding regression faults in applications. In practice there are also likely to be faults left over from earlier versions of the program under test, and the use of a change-based criterion actually may make it more difficult to find these faults.

An additional threat to validity arises from the differences between the change sets identified by JDIFF and CHIANTI as mentioned in Section 3.1. To determine how much impact these differences had on the actual results, we also implemented a changed method coverage based on JDIFF. The differences between the Kendall τ coefficients for the two different changed method implementations were generally small, exceeding 0.1 in only one case (*NanoXML* v4/v5) and 0.01 in 5 other cases, suggesting that the different change analysis techniques had little effect on the results.

6. RELATED WORK

Several different coverage criteria for measuring the adequacy of a test suite have been developed (e.g. [17, 19]), including statement coverage, branch coverage, and various notions of dataflow coverage. In this work, we derive three new criteria, changed statement coverage, changed branch coverage, and changed method coverage, from these traditional coverage criteria, and compare these new criteria to their non-change-based counterpoints.

A number of techniques have been proposed for evaluating different coverage criteria. For example, Weyuker developed a set of 11 properties that a coverage criteria should possess [30] that were later applied to evaluating object-oriented coverage criteria by Perry and Kaiser [18]. Both Wong et al. [33] and Namin and Andrews [16] present studies comparing the effect of test suite size and coverage on fault detection that are similar in setup to ours. Wong et al. find that block coverage correlates better with fault detection than test suite size. Namin and Andrews found that fault detection effectiveness correlated well with block coverage, decision coverage, and two different dataflow coverage criteria; however a model that also included test suite size was even more effective. Unlike our work, these studies are not looking specifically at the practice of regression testing where it is expected that the faults are more likely to be in the changed portions of the code.

Regression testing focuses on ensuring that modifications do not impact pre-existing functionality. Most of the regression testing literature falls into one of two categories: regression test selection (surveyed in [22]) and regression test prioritization (e.g. [9, 11, 12, 13, 24, 27, 29]). Regression test selection attempts to select

a subset of the entire test suite of a program that will identify any faults introduced by the modification into the pre-existing functionalities of the system. Regression test prioritization uses heuristics to reorder the tests in a test suite to increase the likelihood that any newly introduced faults will be revealed earlier in the testing process. These techniques focus on ways of more efficiently using an existing test suite to test a modified version of a system while change-based coverage criteria attempt to evaluate the quality of a test suite relative to a modification and to indicate portions of the modified application that may need additional testing.

Bates and Horwitz [3] and Harrold and Rothermel [21] present techniques for identifying the set of changed dataflow and control-flow testing requirements for a modified program. Their techniques support several criteria including statement coverage, branch coverage, and def-use testing. These techniques use program dependence graph or system dependence graph representations for the original and modified programs to identify the changed testing requirements. While not presented as such, these techniques could form the basis for additional change-based coverage criteria.

The goal of the MATRIX [2] and MATRIXRELOADED [26] projects has been to identify testing requirements for changed software based on data- and control-flow chains and on the symbolic state of the program at the beginning and end of these chains. Then, test inputs are generated to cover these requirements. These coverage criteria have been shown to be more effective and sometimes more cost-effective at identifying changed behavior in the modified programs than changed statement coverage or weaker changed dataflow coverages. However, this work has not directly evaluated the ability of their coverage criteria to expose faults or predict the fault exposure capabilities of different test suites.

7. CONCLUSIONS

Code coverage criteria are commonly used to assess the quality of test suites. The basic idea is that a test suite is likely to be effective at revealing faults if it exercises the code where the fault is located. Therefore, increased code coverage is expected to correlate with more revealed faults. However, achieving full coverage according to traditional coverage criteria is often impossible or impractical when applications contain unreachable code, which may arise due to programmer carelessness, or when the implementation of a feature is incomplete.

In this paper, we start from the popular assumption that a disproportionate number of faults is likely to reside in recently changed code. Based on this assumption, we propose several *change-based coverage criteria* that reflect to what extent changes with respect to a previous program version are exercised by a test suite. In a set of experiments on programs from the SIR repository, we found change-based criteria to reveal faults better than traditional criteria, and to enable the construction of much smaller test suites with similar fault detection effectiveness than those constructed according to traditional coverage criteria. We also reported on a case study that shows that achieving (near) 100% coverage according to a change-based criterion is both feasible and useful.

8. REFERENCES

- [1] T. Apiwattanapong, A. Orso, and M. J. Harrold. JDIFF: A differencing technique and tool for object-oriented programs. *Automated Software Engineering: An International Journal*, 14(1):3–36, 2007.
- [2] T. Apiwattanapong, R. Santelices, P. K. Chittimalli, A. Orso, and M. J. Harrold. MATRIX: Maintenance-oriented testing requirements identifier and examiner. In *Proceedings of*

- Testing: Academic and Industrial Conference - Practice And Research Techniques*, pages 137–146, 2006.
- [3] S. Bates and S. Horwitz. Incremental program testing using program dependence graphs. In *Proceedings of the Symposium on Principles of Programming Languages*, pages 384–396, 1993.
- [4] S. Berner, R. Weber, and R. K. Keller. Enhancing software testing by judicious use of code coverage information. In *Proceedings of the International Conference on Software Engineering*, pages 612–620, 2007.
- [5] J. Bloch. *Effective Java (2nd Edition)*. Prentice Hall PTR, 2008.
- [6] S. A. Bohner and R. S. Arnold. An introduction to software change impact analysis. In S. A. Bohner and R. S. Arnold, editors, *Software Change Impact Analysis*, pages 1–26. IEEE Computer Society Press, 1996.
- [7] H. Do, S. G. Elbaum, and G. Rothermel. Supporting controlled experimentation with testing techniques: An infrastructure and its potential impact. *Empirical Software Engineering: An International Journal*, 10(4):405–435, 2005.
- [8] M. Doliner. COBERTURA.
<http://cobertura.sourceforge.net/>.
- [9] D. Jeffrey and N. Gupta. Test case prioritization using relevant slices. In *Proceedings of the International Computer and Software Applications Conference*, pages 411–420, 2006.
- [10] M. G. Kendall. A new measure of rank correlation. *Biometrika*, 30(1/2):81–93, 1938.
- [11] J. Kim and A. Porter. A history-based test prioritization technique for regression testing in resource constrained environments. In *Proceedings of the International Conference on Software Engineering*, pages 119–129, 2002.
- [12] D. Leon and A. Podgurski. A comparison of coverage-based and distribution-based techniques for filtering and prioritizing test cases. In *Proceedings of the International Symposium on Software Reliability Engineering*, pages 442–453, 2003.
- [13] Z. Li, M. Harman, and R. M. Hierons. Search algorithms for regression test case prioritization. *IEEE Transactions of Software Engineering*, 33(4):225–237, 2007.
- [14] J. C. Miller and C. J. Maloney. Systematic mistake analysis of digital computer programs. *Communications of the ACM*, 6(2):58–63, 1963.
- [15] G. J. Myers. *The Art of Software Testing*. John Wiley and Sons, Inc, 1979.
- [16] A. S. Namin and J. H. Andrews. The influence of size and coverage on test suite effectiveness. In *Proceedings of the International Symposium on Software Testing and Analysis*, pages 57–68, 2009.
- [17] S. C. Ntafos. On required element testing. *IEEE Transactions on Software Engineering*, 10(6):795–803, 1984.
- [18] D. E. Perry and G. E. Kaiser. Adequate testing and object-oriented programming. *Journal of Object Oriented Programming*, 2(5):13–19, 1990.
- [19] S. Rapps and E. J. Weyuker. Selecting software test data using data flow information. *IEEE Transactions on Software Engineering*, 11:367–375, 1985.
- [20] X. Ren, F. Shah, F. Tip, B. G. Ryder, and O. Chesley. CHIANTI: A tool for practical change impact analysis of Java programs. In *Proceeding of the Conference on Object Oriented Programming, Systems, Languages and Applications*, pages pp 432–448, 2004.
- [21] G. Rothermel and M. J. Harrold. Selecting tests and identifying test coverage requirements for modified software. In *Proceedings of the International Symposium on Software Testing and Analysis*, pages 169–184, 1994.
- [22] G. Rothermel and M. J. Harrold. Analyzing regression test selection techniques. *IEEE Transactions on Software Engineering*, 22:529–551, 1996.
- [23] G. Rothermel and M. J. Harrold. A safe, efficient regression test selection technique. *ACM Transactions on Software Engineering and Methodology*, 6(2):173–210, 1997.
- [24] G. Rothermel, R. Untch, C. Chu, and M. J. Harrold. Test case prioritization. *IEEE Transactions of Software Engineering*, 27(10):929–948, 2001.
- [25] B. G. Ryder and F. Tip. Change impact analysis for object-oriented programs. In *Proceeding of the Workshop on Program Analysis for Software Tools and Engineering*, pages 46–53, 2001.
- [26] R. Santelices, P. K. Chittimalli, T. Apiwattanapong, A. Orso, and M. J. Harrold. Test-suite augmentation for evolving software. In *Proceedings of the International Conference on Automated Software Engineering*, pages 218–227, Sept. 2008.
- [27] A. Srivastava and J. Thiagarajan. Effectively prioritizing tests in development environment. In *Proceedings of the International Symposium on Software Testing and Analysis*, pages 97–106, 2002.
- [28] M. Stoerzer, B. G. Ryder, X. Ren, and F. Tip. Finding Failure-inducing Changes in Java Programs Using Change Classification. In *Proceedings of the Symposium on the Foundations of Software Engineering*, pages 57–68, 2006.
- [29] A. Walcott, M. L. Soffa, G. M. Kapfhammer, and R. Roos. Time-aware test suite prioritization. In *Proceedings of the International Symposium on Software Testing and Analysis*, pages 1–12, 2006.
- [30] E. J. Weyuker. The evaluation of program-based software test data adequacy criteria. *Communications of the ACM*, 31(6):668–675, 1988.
- [31] J. Wloka, B. G. Ryder, and F. Tip. JUNITMX - a change-aware unit testing tool. In *Proceedings of the International Conference on Software Engineering*, pages 567–570, 2009.
- [32] J. Wloka, B. G. Ryder, F. Tip, and X. Ren. Safe-commit analysis to facilitate team software development. In *Proceedings of the International Conference on Software Engineering*, pages 507–517, 2009.
- [33] W. E. Wong, J. R. Horgan, S. London, and A. P. Mathur. Effect of test set size and block coverage on the fault detection effectiveness. In *Proceedings of the International Symposium on Software Reliability Engineering*, pages 230–238, 1994.
- [34] Q. Yang, J. J. Li, and D. M. Weiss. A survey of coverage based testing tools. In *Proceedings of the International Workshop on Automation of Software Test*, pages 99–103, 2006.