

Reducing the Costs of Bounded-Exhaustive Testing

Vilas Jagannath, Yun Young Lee, Brett Daniel, and Darko Marinov

Department of Computer Science, University of Illinois at Urbana-Champaign
Urbana, IL 61801, USA

{vbanga12,lee467,bdaniel13,marinov}@cs.uiuc.edu

Abstract. Bounded-exhaustive testing is an automated testing methodology that checks the code under test for *all inputs* within given bounds: first the user describes a set of test inputs and provides test oracles that check test outputs; then the tool generates all the inputs, executes them on the code under test, and checks the outputs; and finally the user inspects failing tests to submit bug reports. The costs of bounded-exhaustive testing include machine time for test generation and execution (which translates into human time waiting for these results) and human time for inspection of results. This paper proposes three techniques that reduce these costs. Sparse Test Generation skips some tests to reduce the time to the first failing test. Structural Test Merging generates a smaller number of larger test inputs (rather than a larger number of smaller test inputs) to reduce test generation and execution time. Oracle-based Test Clustering groups failing tests to reduce the inspection time. Results obtained from the bounded-exhaustive testing of the Eclipse refactoring engine show that these three techniques can substantially reduce the costs while mostly preserving fault-detection capability.

1 Introduction

Testing is an important but expensive part of software development, estimated to take more than half of the total development cost [1]. One approach to reducing the cost is to automate testing. Bounded-exhaustive testing is an automated approach that checks the code under test for *all inputs* within given bounds [2, 3, 4, 5, 6]. The rationale is that many faults can be revealed within small bounds [7, 8], and exhaustively testing within the bounds ensures that no “corner case” is missed. Bounded-exhaustive testing has been used in both academia and industry to test several real-world applications, with some recent examples including testing of refactoring engines [5] and a web-traversal code [6].

Bounded-exhaustive testing consists of three activities. First, the user describes a set of test inputs and provides test oracles that check test outputs. Second, the tool generates all the inputs, executes them on the code under test, and checks the outputs using the oracles. Third, the user inspects failing tests to submit bug reports or debug the code; typically, bounded-exhaustive testing produces a large number of *failures* for each *fault* found. Two key costs in this

context are machine time for test generation and execution (which also translates into human time for waiting for these results [9,10]) and human time for inspection of failures. Previous experience shows that bounded-exhaustive testing can discover important faults [3,4,5,11] but also can have high costs.

This paper proposes, and evaluates on a case study, three novel techniques that reduce these costs of bounded-exhaustive testing. The three techniques address various costs and can be used individually or synergistically.

Sparse Test Generation (STG): We present a new technique that reduces the time to first failure (abbreviated *TTF*), i.e., the time that the user has to wait after starting a tool for bounded-exhaustive testing until the tool finds *a* failing test. Note that in this context there is usually a large number of failing tests (say, hundreds or even thousands) or no failing test (if the code under test reveals no fault for any generated test). TTF measures only the time to the *first* failure (not all failures). It is an important practical metric that captures the user idle time. Previous research shows, in a related context of regression testing, that reducing the time to failure can significantly help in development [9,10]. STG works by making two passes through test generation. The first, *sparse*, pass skips some tests in an attempt to reduce TTF. While this pass is related to test suite minimization/reduction/prioritization [12,13,14,15,16,17], the main challenge is to skip tests while they are being generated and not to select some tests only after all have been generated. The second, *exhaustive*, pass generates all the tests to ensure exhaustive checking within the given bound. Effectively, STG trades off (substantially) decreasing TTF for (slightly) increasing the total time.

Structural Test Merging (STM): We present a new technique that reduces the total time for test generation and execution. In bounded-exhaustive testing, users typically describe a test set with a *large number of small tests*, while we advocate considering test sets with a *smaller number of larger tests*. Our technique is inspired by the work on test granularity [18,19] which studied the cost-benefit trade-offs in using a larger number of smaller tests versus a smaller number of larger tests. That work mostly considered manually written tests for regression testing, while we focus on automatically generated tests. Moreover, that work considered cases where larger tests can be automatically built from smaller tests by simply *appending* (e.g., if each test is a sequence of commands, a longer test sequence can be obtained by simply appending a number of shorter test sequences), while we consider cases where it is harder to build larger tests from smaller tests (e.g., simply appending two test input programs together while testing a compiler or a refactoring engine would likely result in a compilation error as these programs have program entities with the same name; moreover, renaming would reduce the opportunity of speeding up test execution). Instead of simply appending tests, our technique *merges* them based on their structure, hence the name STM.

Oracle-based Test Clustering (OTC): We present a new technique that reduces the human time for inspection of failing tests. Bounded-exhaustive testing can produce a large number of failing tests, and a tester/developer has to map these failures to distinct faults to submit bug reports or debug the code under

test. Our technique builds on the ideas from test clustering [20, 21, 22, 23, 24, 25] where the goal is to split (failing) tests into groups such that all tests in the same group are likely due to the same underlying fault. Previous work mostly considered manually written tests or actual programs runs, and clustering was based on *execution profiles* obtained from monitoring test execution. In contrast, we consider automatically generated test inputs, and our technique exploits information from oracles. Typically, an oracle only states *if* some test passed or failed, i.e., the output from an oracle is a boolean. However, in some domains oracles also state *how* the result is incorrect, i.e., the output from an oracle is an error *message*. OTC splits tests based on oracle messages, and our results suggest that it is beneficial to build such oracles whenever possible. The key to our technique is *abstracting* messages and not comparing them directly.

Case Study: We implemented our three new techniques in the ASTGen framework for bounded-exhaustive testing of refactoring engines [5]. We chose ASTGen for three reasons: it had enabled finding actual faults in real large software (we had found a few dozens of new faults in the refactoring engines of Eclipse and NetBeans, two popular IDEs for Java [5]); we were familiar with the framework; and we personally experienced the costs of using ASTGen. We evaluated the techniques on testing 6 refactorings with 9 generators (explained later in the text). The results show that (1) STG can reduce TTFB almost 10x (an order of magnitude) when there is a failure, while increasing the total test generation and execution time only 10% when there is no failure; (2) STM can reduce the total time 2x-6x (in one instance from over 6 hours to 70 minutes) and even more (but with some reduction of the fault-detection capability); and (3) OTC can reduce the number of tests to be inspected by clustering hundreds of failing tests into a few groups (up to 11) such that almost all tests within the group are due to the same fault. In summary, the results show that the new techniques can substantially reduce both the machine time and the human time without reducing the fault-detection capability.

2 Example

To illustrate how our techniques reduce the costs of bounded-exhaustive testing, we discuss testing of the PullUpMethod refactoring in the Eclipse refactoring engine using the ASTGen framework. We first describe what PullUpMethod is. We then describe how to use ASTGen for bounded-exhaustive testing of this refactoring. We finally discuss how our new techniques improve on ASTGen.

Each refactoring is a program transformation that changes the program code but not its external behavior [26]. Programmers undertake refactorings to improve design of their programs. For example, PullUpMethod is a refactoring that moves a method from some class into one of its superclasses (usually because the same method is useful for other subclasses of that superclass). Figure 1 shows a simple application of the PullUpMethod refactoring. Note that moving the method also requires properly updating the references within the method body, i.e., replacing `super.f` with `this.f`.

<pre>// Before refactoring class A { int f; } class B extends A { void m() { super.f = 0; } }</pre>	<pre>// After refactoring class A { int f; void m() { this.f = 0; } } class B extends A { }</pre>	<pre>// Before refactoring class A { } class B extends A { int f; void m() { this.f = 0; } }</pre>	<pre>// Refactoring engine // warning: // Cannot pull up: // method 'm' // without pulling up: // field 'f'</pre>
---	---	--	---

Fig. 1: Example applications of the PullUpMethod refactoring

Refactoring engines are development tools that automate applications of refactorings. They are an important part of modern IDEs such as Eclipse [27]. To apply PullUpMethod, the developer instructs the engine which method to move to which superclass in the *input* program. The engine first checks whether the move is permitted (e.g., PullUpMethod should not move a method to a superclass if the superclass already has a method with the same signature). If it is, the engine appropriately transforms the program. The *output* is either a transformed program or a set of warning messages that indicate why the move would not be permitted, as illustrated in Figure 1.

Testing the implementation of PullUpMethod requires generating a number of input programs, invoking the refactoring engine on them, and checking whether it gives the appropriate output (either a correctly transformed program or an expected set of warning messages). Testers can have good intuition about which input programs could reveal a fault. For instance, PullUpMethod may have faults if the subclass and superclass have some additional relationship, e.g., being an inner or a local class or being related through a third class. Also, there may be faults for some expressions and statements that include field and method references from the body of the method being pulled up or to the method being pulled up. However, it is time-consuming and error-prone to manually generate a large number of such input programs.

We previously developed the ASTGen framework for bounded-exhaustive testing of refactoring engines [5]. ASTGen allows the tester to write *generators* that can automatically produce a (large) number of (small) input programs for testing refactorings. ASTGen generates *all* these inputs, executes the refactoring engine on them, runs several oracles to validate the outputs, and reports failures.

For instance, to test PullUpMethod, we can use a generator that produces programs with three classes in various relationships. For this specific case, ASTGen generates 1,152 input programs, of which 160 result in failing oracles. A detailed inspection of these failures shows that they reveal 2 distinct faults. While finding these faults is clearly positive, there are costs. Test generation and execution (including oracles) take about 27 minutes (on a typical desktop), and the time to find the first failure is about 9 minutes. Also, identifying the 2 distinct faults among 160 failing tests is labor-intensive and tedious.

This paper proposes three techniques that reduce these costs. *STG* addresses the time to first failure (TTFF) by first sampling some inputs rather than exhaustively generating all inputs from the beginning. For our specific example,

TTFB is on average reduced almost an order of magnitude, from about 9 minutes to 1 minute. *STM* addresses the total time for test generation and execution. Instead of testing PullUpMethod for 1,152 (small) programs that exercise various features in isolation, STM builds larger programs that combine some of the features, e.g., combine several expressions or statements that include field and method references to/from the method being pulled up. The tester can choose how many features to combine. In this example, the least aggressive combination reduces the total time from 27 minutes to about 4 minutes, and the most aggressive combination reduces the total time further to under 1 minute. *OTC* addresses the cost of failure inspection. It clusters the failing tests into groups that are likely to be due to the same fault, and thus the tester can inspect only one or a few tests from these “equivalence classes”. Our clustering is based on oracle messages and can consider more or fewer details of the messages. The basic clustering splits 160 failing tests into 127 clusters, but our best clustering splits them into just 3 clusters that reliably find the 2 faults. In contrast, random sampling could miss faults, e.g., one of our experiments shows that it finds on average 1.77 out of 2 faults in this case.

3 Background: ASTGen

We now describe in more detail two parts of the ASTGen framework that are relevant to present the three techniques introduced in this paper. ASTGen allows the testers to write *generators*—pieces of code that implement a specific interface—which ASTGen runs to automatically generate input programs. ASTGen then applies refactorings on these inputs and runs the *oracles* on the outputs.

Generators: Each generator is a piece of Java code that produces elements of Java abstract syntax trees (ASTs), which can be pretty-printed as Java source. Conceptually, generators are close to grammar-based generation [28, 29], but ASTGen uses Java code rather than a grammar formalism as explained elsewhere [5]. ASTGen provides (1) a large library of basic generators, (2) several mechanisms to compose and link simpler generators into more complex generators, and (3) customization of generators using Java code. Some of the generators that ASTGen provides include:

Field Declaration Generator produces many different field declarations that vary in terms of type (`int`, `byte`, `boolean`, array or non array, etc.), visibility (`private`, `public`, etc.), and name of the declared field.

Field Reference Expression Generator is linked to the Field Declaration Generator and produces different expressions that reference the declared field in various ways, including field accesses and operations (`this.f`, `new A().f`, `super.f`, `f++`, `!f`, etc.).

Single Class Field Reference Generator is composed on top of the Method Declaration Generator and produces classes with one field (obtained from the Field Declaration Generator) and one method that references the field in various ways.

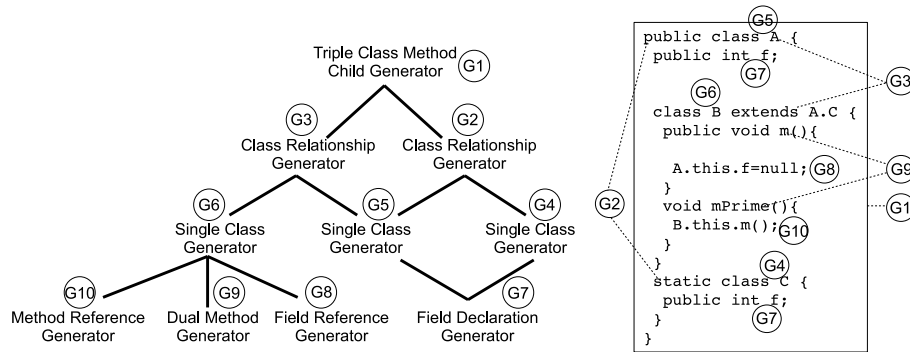


Fig. 2: Triple Class Method Child Generator structure and a generated test input

Dual Class Relationship Generator is composed upon generators that produce classes (e.g., Single Class Field Reference Generator) and produces two classes with various relationships between them (inheritance, inner class, local class, etc.).

While the main purpose of generators is to actually produce the test inputs, they also encode the space of all inputs to be produced. Consider this scenario:

Inputs: Programs with three classes A, B, and C.

- B extends C; B has a method m and a method mPrime that invokes m.
- C and A each have a field f that may be referenced by m.

Test: Pull up method m from class B to class C.

The user can generate all these inputs by writing a generator that composes and links several library generators. Figure 2 shows the overall structure of a generator, called Triple Class Method Child Generator, that encodes this input space. The figure also shows a sample test input produced by this generator and how the input sub-parts match the sub-generators responsible for producing them. By iterating through all the variations of the sub-generators, the Triple Class Method Child Generator produces 1,152 test inputs.

Oracles: While generators are the core of ASTGen and help testers to produce a large number of input programs for testing refactorings, it would be impractical that the testers check the result of each refactoring application. Oracles automate checking of the results so that the testers only have to inspect a smaller number of tests that fail the oracles (and likely detect faults). ASTGen provides two generic oracles and allows the users to write refactoring-specific oracles:

Compilation Failure Oracle flags tests where the refactored program has a compilation error: if the input program compiles, then the output program should also compile.

Erroneous Warning Oracle flags tests where the refactoring engine raised a warning about a refactoring application, but ignoring that warning results in a refactored program with no compilation errors (or custom failures).

Custom Oracles are specific to the refactoring being applied. For example, a custom oracle for `RenameMethod` checks that renaming a method, say `m` to `p`, and then renaming back, `p` to `m`, results in the same program.

The output of traditional oracles are only booleans (pass or fail), but the AST-Gen oracles can provide additional information about the failure, e.g., messages from the compiler or warnings from the refactoring engine.

4 Sparse Test Generation (STG)

Generators can encode and produce all the test inputs within defined bounds. Bounded-exhaustive testing checks the code under test for all these inputs. This usually consumes a large amount of machine time since the number of inputs generated is fairly large. For example, ASTGen generators can generate thousands of test inputs, and it can take hours of machine time to execute the refactorings on all those inputs. Additionally, this time translates into human time required by the developer to wait for the execution of the tests to complete.

Note that as soon as a tool reports a failure, the developer can start inspecting it to file a bug report or to debug the fault that caused the failure. In theory, the time the tool takes for generation and testing after the first failure is not important since the developer does not have to idle. For this reason, we consider the Time to First Failure (TTFF) as the key metric in interactive bounded-exhaustive testing. If no generated test input results in a failure, the developer has to wait for the entire generation and testing process to complete.

STG is our technique that aims to reduce the TTFF. STG has two phases: **Sparse Generation** is motivated by our observation that failing test inputs are often located closely together in the sequence of inputs produced by a generator, and thus, to find a failure, it is often not necessary to exhaustively generate all the inputs but only one input from a closely located group. Therefore, this phase makes “jumps” through the generation sequence. The jump length is not constant (since the failing tests may be in a stride that a constant jump would miss) but *each jump is (uniformly) random within some length limit*. The key is to determine an appropriate limit: a lower limit increases the overhead of STG compared to the basic, *dense* bounded-exhaustive testing, while a higher limit decreases the chance that Sparse Generation finds a failure (and thus increases the TTFF). We use the limit of 20 as it provides a good trade-off: the expected jump is of length $(1+20)/2$, which increases the total time by less than 10% when there is no failure. If Sparse Generation finds a failing test, it usually does so quickly; the results from Section 7 show that STG reduces the TTFF by an order of magnitude in most cases compared to the dense generation. However, STG is a heuristic and, in general, could keep missing failures until the very end while dense generation would have found those failures at the very beginning. **Exhaustive Generation** follows Sparse Generation and does basic bounded-exhaustive testing (1) to ensure that a failing test input will be found if one exists and (2) to find all the failing tests that Sparse Generation missed (which can help in clustering failures or debugging [20, 23, 25]).

<pre> public class A { public int f; class B extends C { private void m(){ this.f=0; } void mPrime(){ m(); } } } class C { public int f; } </pre>	<pre> public class A { public int f; class B extends C { private void m(){ new A().f=0; } void mPrime(){ m(); } } } class C { public int f; } </pre>	<pre> public class A { public int f; class B extends C { private void m(){ super.f=0; } void mPrime(){ m(); } } } class C { public int f; } </pre>
---	--	--

Fig. 3: Unmerged test inputs

5 Structural Test Merging (STM)

TTFB is an important metric in bounded-exhaustive testing. Another important metric is the total time for test generation and execution. This time can be very long when generators produce a large number of inputs, which is the case for typical top-level ASTGen generators. For example, consider the number of inputs for the Triple Class Method Child Generator shown in Figure 2. Each of its sub-generators has a small number of variations—G2 has 2 (inner, outer); G3 has 3 (inner, method inner, outer); G4, G5, and G6 have 1; G7 has 2 (public, private); G8 has 6 (f, new A().f, A.this.f, etc.); G9 has 4 (public, private, same/different signature); and G10 has 4 (m(), new B().m(), this.m(), B.this.m())—but the top-level generator produces $2 \times 3 \times 1 \times 1 \times 1 \times 2 \times 6 \times 4 \times 4 = 1152$ combinations.

STM reduces the number of test inputs while still aiming to preserve their exhaustiveness: instead of producing a large number of small input programs, STM produces a smaller number of larger input programs by *merging* appropriate program elements. For example, the Triple Class Method Child Generator produces the three inputs shown in Figure 3. The only difference between the three are the highlighted statements, generated by the Field Reference sub-generator (G8). Figure 4 shows an input that contains all these three statements. This single, merged input encodes the same input space as the three unmerged inputs. This structural merging transformation is the crux of our STM technique.

STM exploits the compositional structure of the sub-generators to produce merged test inputs. Figure 4 shows an alternative structure for the Triple Class Method Child Generator: a new Field Reference Merging Single Class Generator (G8M) merges together all the program elements produced by the original Field Reference Generator (G8). While figures 2 and 4 show a generator before and after a *single* application of the structural merging transformation, it is possible to apply the transformation *multiple* times within the hierarchical structure of a generator. Each application leads to a multiplicative reduction in the number of generated inputs. For example, the original Method Reference Generator (G10) can also be modified to a generator G10M that merges together all the different method invocation statements. Together, G8M and G10M produce inputs that merge *both* field references and method references. We refer to the number of

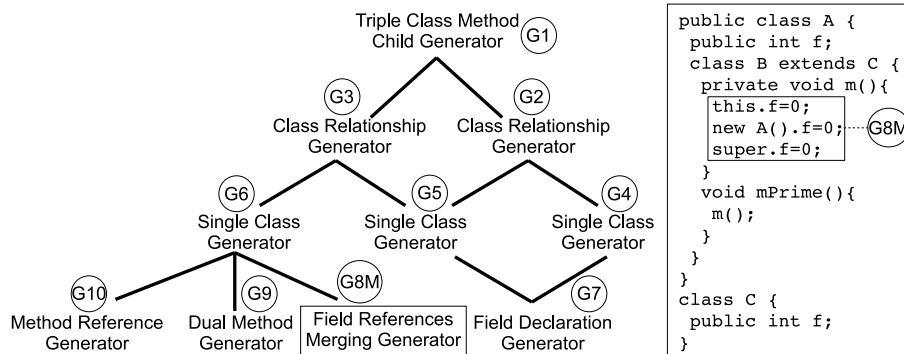


Fig. 4: Merged generator structure and a generated merged test input

transformation applications as *merging level*: for the Triple Class Method Child Generator, merging level M1 has only G8M, and merging level M2 has both G8M and G10M. The unmerged generator produces 1,152 inputs, and levels M1 and M2 reduce the number of inputs to 192 and 48, respectively.

While STM achieves significant time savings, it is important to note its two potential drawbacks. One potential drawback is that larger inputs, through the interference of program elements, can mask some test failures [18, 19]. Consider, for example, merging together all the different field references (as in Figure 4). There may be a failure triggered by one of the field reference statements which gets masked by the presence of the other field reference statements. However, this interference can also go the other way: larger inputs may trigger new failures that smaller inputs do not trigger. The other drawback is the effect of larger inputs on debugging. STM produces fewer larger inputs rather than more smaller inputs, but (failing) smaller inputs typically make it easier to perform fault localization. We could take two approaches to address this. One approach is to reduce inputs by applying Delta Debugging [30] on the larger failing input to try to isolate the part of the input that triggers the failure. Another approach, enabled by the fact that larger inputs are produced by merging generators, is to regenerate the small inputs that represent the larger failing input.

6 Oracle-based Test Clustering (OTC)

The experience with bounded-exhaustive testing in academia and industry shows that it can find faults in real code [4, 3, 5, 11] but also produces a large number of failures. Identifying a few faults out of many failures is a challenging task. OTC is a new technique that helps in this task by splitting failing tests into groups such that all tests in the same group are likely due to the same fault.

OTC exploits information from oracles. Recall that ASTGen oracles provide messages about the failures, e.g., if a refactored program does not compile, ASTGen reports the compilation error provided by the compiler. We use these

messages to cluster the failing tests by grouping together those tests that have *exactly the same messages*. (A test can produce multiple messages, which our experiments compare as lists, not bags or sets.) However, directly using *concrete* messages provided by the compiler can result in a large number of small clusters, e.g., two compilation errors may differ only in line or column numbers, say, “3:8:field f not visible” and “2:6:field f not visible”. Instead, we use *abstract* messages that ignore some details such as line and column numbers. One can further consider ignoring exact messages and clustering based on *which* oracle failed, *not how* it failed. The trade-off is that creating too many clusters increases inspection effort, while creating too few clusters increases the chance to miss a fault. Our evaluation compares four clustering options: Concrete Message, Abstract Message, Oracle Name, and Random Selection (a base case with no clustering).

7 Case Study

We evaluated our three new techniques in the ASTGen framework for bounded-exhaustive testing of refactorings engines. We tested 6 refactorings using 9 generators listed in Figure 5. For each generator and several merging levels, we tabulate the number of inputs generated, various times and APFD metric (described below), the number of failing inputs, and the number of distinct faults. We previously tested these refactorings with these generators and found a number of faults [5]. The goal of this study was to evaluate whether the new techniques reduce the testing costs, but due to OTC, we also found a new fault in the PushDownMethod refactoring, previously missed [5] due to random sampling. We ran all experiments in Eclipse 3.3.2 on a dual core 3.4GHz machine.

Sparse Test Generation (STG): Figure 5 shows the time results for ASTGen with and without STG. The ‘Dense’ subcolumns show the total time and time to first failure (TTFF) for bounded-exhaustive testing without STG. If no failure exists, TTFF shows ‘n/a’. The ‘Sparse’ column shows average values for TTFF if a failure exists (roughly the top half of the table) and the total time if no failure exists (the bottom half of the table). These times are averaged over 20 random seeds, with the jump limit of 20, as discussed in Section 4. The main questions about STG are how it affects TTFF and the total time.

STG reduces TTFF in all cases where the dense TTFF was significant (a minute or more): the speedup ranges from 9.00x to 10.58x, with an average of an order of magnitude. In a few cases with very small dense TTFF, STG had a slowdown of at most 1 sec. Recall the two phases of STG; the reduction in TTFF implies that the sparse phase found a failure before the exhaustive phase.

STG increases the total time, as expected. With the jump limit of 20, the overhead of the additional sparse phase is expected to be slightly under 10% of the total time for dense generation. Our experiments confirm that this is indeed the case: the slowdown ranges from -6.48% to -9.60%. In summary, STG achieves a 10x speedup in TTFF for only a 10% slowdown in the total time.

We further evaluated STG using the Average Percentage Fault Detected (APFD) metric introduced by Rothermel et al. [12] to compare techniques for

Refactoring	Generator	ML	Inputs	Dense		Sparse	APFD [%]		Failures	Faults
				Time	TTF		Dense	Sparse		
PushDown-Field	DualClass-FieldReference	M0	7416	133:32	7:33	0:47	74.98	98.16	1074	2
		M1	1236	22:43	0:01	0:02	99.23	88.62	179	1
		M2	12	0:21	0:00	0:01	95.83	74.17	4	1
		M3	3	0:13	0:00	0:00	83.33	63.33	1	1
Encapsulate-Field	DualClass-FieldReference	M0	23760	427:09	73:34	7:14	58.03	97.59	486	3
		M1	3960	71:50	12:03	1:11	69.82	97.77	354	3
		M2	72	1:19	0:13	0:03	74.31	80.56	31	2
		M3	18	0:26	0:06	0:03	58.33	73.15	8	2
	SingleClass-FieldReference	M0	8576	155:15	0:22	0:03	75.37	97.61	836	4
		M1	2144	39:04	0:21	0:03	66.86	97.59	242	4
		M2	1072	19:35	0:09	0:02	84.25	93.04	144	3
		M3	268	4:55	0:02	0:02	72.70	88.11	62	3
PushDown-Method	DualClass-MethodParent	M0	960	22:19	11:28	1:05	43.91	93.59	180	3
		M1	192	4:07	2:07	0:14	41.75	91.89	38	3
		M2	48	0:45	0:28	0:21	40.63	87.85	2	1
PullUp-Method	TripleClass-MethodChild	M0	1152	27:02	9:09	1:01	13.19	95.77	160	2
		M1	192	3:57	1:25	0:09	48.18	95.36	96	2
		M2	48	0:47	0:17	0:02	56.25	89.58	24	2
	DualClass-MethodChild	M0	576	13:22	n/a	14:14	n/a	n/a	0	0
		M1	96	1:49	n/a	1:55	n/a	n/a	0	0
		M2	24	0:21	n/a	0:22	n/a	n/a	0	0
Rename-Field	DualClass-FieldReference	M0	23760	629:01	n/a	689:17	n/a	n/a	0	0
		M1	3960	107:26	n/a	117:48	n/a	n/a	0	0
		M2	72	1:56	n/a	2:04	n/a	n/a	0	0
		M3	18	0:34	n/a	0:34	n/a	n/a	0	0
	SingleClass-FieldReference	M0	8576	229:00	n/a	250:59	n/a	n/a	0	0
		M1	2144	57:28	n/a	62:56	n/a	n/a	0	0
		M2	1072	28:44	n/a	31:28	n/a	n/a	0	0
		M3	268	7:15	n/a	7:57	n/a	n/a	0	0
		M0	9540	173:32	n/a	190:11	n/a	n/a	0	0
Rename-Method	SingleClass-MethodReference	M1	4900	89:26	n/a	98:05	n/a	n/a	0	0
		M2	140	2:37	n/a	2:50	n/a	n/a	0	0
		M3	80	1:31	n/a	1:37	n/a	n/a	0	0

Fig. 5: Sparse Test Generation and Structural Test Merging Results
Legend: ML = Merging Level, TTF = Time to First Failure, All times in minutes:seconds

test prioritization and extended by Walcott et al. [16] for test selection. APFD measures the number of faults detected in terms of the number of tests executed, whereas TTF is based on the first failure (not all faults) and actual time (not number of tests) as TTF aims to capture the waiting time for testers in interactive bounded-exhaustive testing, similar to recent extensions of APFD [10]. APFD ranges between 0 and 100%, with higher values being better. Figure 5 shows APFD, with ‘Sparse’ averaged over 20 random seeds. The results show that STG improves APFD in all cases where the dense TTF was significant.

Structural Test Merging (STM): Figure 5 shows the results for STM for several merging levels of each of the generators. The merging level number (e.g., 3 in M3) represents the number of structural merging transformations applied to the unmerged generator (labeled M0) to obtain the corresponding merged generator, as discussed in Section 5. The main questions about STM are how it affects times (total and TTF) and the number of failures/faults detected.

Each merging level reduced both the total time and TTF compared to its previous level and thus to M0. On average, level M1 achieved 5x speedup, and level M2 achieved 130x speedup compared to M0 for the total time. The merged

Refactoring	Generator	ML	Random		Oracle		Abstract		Concrete	
			FD	NC	FD	NC	FD	NC	FD	NC
PushDownField	DualClassFieldReference	M0	1.99	1	2	2	2	5	2	68
		M1	1	1	1	1	1	3	1	59
		M2	1	1	1	1	1	2	1	4
		M3	1	1	1	1	1	1	1	1
EncapsulateField	DualClassFieldReference	M0	2.24	1	2.24	2	3	4	3	51
		M1	2.05	1	2.31	2	3	4	3	112
		M2	1.73	1	2	2	2	4	2	8
		M3	1.75	1	2	2	2	3	2	3
	SingleClassFieldReference	M0	2.84	1	3.11	2	4	4	4	71
		M1	2.48	1	2.46	2	4	4	4	73
		M2	2.26	1	2.26	1	3	4	3	58
M3		2.30	1	2.30	1	3	5	3	24	
PullUpMethod	TripleClassMethodChild	M0	1.77	1	1.77	1	2	3	2	127
		M1	1.56	1	1.56	1	2	2	2	84
		M2	1.62	1	1.62	1	2	2	2	24
		M3	1.62	1	1.62	1	2	2	2	24
PushDownMethod	DualClassMethodParent	M0	2.19	1	3	2	3	11	3	20
		M1	2.56	1	2.54	2	3	10	3	16
		M2	1	1	1	1	1	1	1	2

Fig. 6: Oracle-based Test Clustering Results
Legend: ML = Merging Level, FD = Faults Detected, NC = Number of Clusters

generators also substantially reduced the TTFF: on average, level M1 achieved 80x speedup, and level M2 generators achieved 150x speedup compared to M0.

Merging did not expose any new faults, but aggressive merging did mask some faults. In particular, level M1 masks only one fault (in PushDownField), but levels M2 and higher mask a much larger number of faults. However, even the highest level of merging finds at least one fault (when there is a fault at M0). Additionally, if one considers TTFF as the most important metric, masking faults at the higher merging levels is not detrimental but actually beneficial: the user can start the exploration from a high level, quickly find failures, and start inspecting them, while the tool continues the exploration at a lower level. In summary, STM can substantially improve total time and TTFF while somewhat reducing the fault-detection capability of bounded-exhaustive testing.

Oracle-based Test Clustering (OTC): Figure 6 shows the results for the four clustering options discussed in Section 6. For each option, we present the number of clusters formed and distinct faults detected by inspecting a number of randomly selected tests from each cluster. The results are averaged over 1000 random seeds. For this experiment, we needed to choose a sampling strategy [20], which determines how many tests to select and from which clusters. The basic strategy, one-per-cluster, selects one test for each cluster; we used this strategy for Abstract Message and Concrete Message. For Random Selection and Oracle Name, which have fewer clusters, we used a strategy that selects more tests per cluster, specifically selects at least as many tests as Abstract Message selects (i.e., the number of clusters that Abstract Message has) and at most 1% of all failing tests. The main questions about OTC are how it affects the number of failures that need to be inspected and the number of faults detected.

To measure the number of faults detected by a set of selected tests, we had to map failing tests to the fault(s) they detect and also had to determine

which faults are *distinct*. We performed two steps. First, a researcher (the second paper author) manually inspected all tests from each cluster (based on Abstract Message) with less than 30 tests and inspected at least 10 tests from each cluster with more than 30 tests. Since all inspected tests from each cluster detected the same fault(s), we extrapolated that all tests in a cluster can detect the same fault(s). We also patched 6 of these faults in Eclipse and confirmed their results from the first step. Second, we asked a researcher (unaware of the details of this study but with a multi-year experience with Eclipse refactorings) to label the faults collected in the first step as potential duplicates of each other or non-faults. This resulted in 12 distinct faults that we used in our experiments.

Abstract Message substantially reduces the number of tests to be inspected to find all the faults, e.g., PullUpMethod for M0 has 160 failing tests, but Abstract Message splits them into 3 clusters, and selecting any 3 tests, one from each cluster, always reveals all 2 faults. The results show that Abstract Message finds *all faults* that Concrete Message finds but requires inspection of much fewer tests, up to over an order of magnitude for lower merging levels. Also, Abstract Message finds *more faults* than Random Selection and Oracle Name while the same number or even fewer tests are inspected. In summary, Abstract Message was the most effective OTC option among the four we compared.

8 Related Work

There is a large body of work on automated testing. Our focus is on bounded-exhaustive testing [2, 3, 4, 5, 6, 11] that tests the code for all inputs within given bounds. Previous work considered how to describe a set of inputs (using declarative [2, 4] or imperative [5] approaches) and how to efficiently generate them. Bounded-exhaustive testing has been successfully used to reveal faults in several real applications [3, 4, 5, 11], but it has costs in machine time for test generation and execution and human time for inspection of failures. This paper presents three new techniques that reduce the costs of such testing.

STG is related to work on test selection/reduction/prioritization [12, 13, 14, 15, 16, 17, 31, 32] whose goal is to reduce the testing cost or to find faults faster by selecting a subset of tests from a test suite and/or ordering these tests. The previous techniques mostly consider regression testing where a test suite exists a priori, and the simplest techniques can randomly select or order these tests. In contrast, STG selects tests while they are being generated, and generation proceeds in a particular order, so arbitrary random sampling is not possible. Finally, STG does not compromise the fault-finding ability [32].

STM is related to work on test granularity [18, 19] which studied the cost-benefit trade-offs in testing with a larger number of smaller tests versus a smaller number of larger tests. The key difference is that previous work considered tests that can be easily *appended* while we consider tests that need to be *merged*. Note that appending tests only saves setup and teardown costs [19], while merging can also reduce test execution cost (e.g., merging 1,152 input programs into 192 input programs requires only 192 applications of the PullUpMethod refactoring).

However, the results are similar in both contexts: larger tests reduce the testing time, but too large tests may miss faults.

OTC is related to work on test clustering/filtering/indexing [20, 21, 22, 23, 24, 25]. Previous work performed clustering based on *execution profiles*, obtained from monitoring test execution. The main novelty of our technique is to exploit information-rich oracles, rather than execution profiles, to cluster failing tests. Our goal is to cluster failing tests to help in identifying the underlying faults. Dickinson et al. [20] present an empirical study that evaluates somewhat different techniques whose goal is to find failures among executions by using cluster analysis of execution profiles. Effectively, those techniques use cluster analysis as approximate oracles. Their results show that cluster filtering of executions can find failures more effectively than random sampling, and that clustering of executions can distinguish failing executions from passing ones.

9 Conclusions

Bounded-exhaustive testing checks the code for all inputs within given bounds. It can find faults but at potentially high costs, including machine time to generate and run tests, and human time to wait for the test results and to inspect failures. We presented three techniques that reduce these costs: Sparse Test Generation skips some tests to reduce the time to first failure by an order of magnitude; Structural Test Merging generates larger tests to reduce test generation and execution time by order(s) of magnitude; and Oracle-based Test Clustering groups failing tests to reduce the inspection time by order(s) of magnitude.

Acknowledgments. We thank Danny Dig for inspecting the faults we found in Eclipse, the anonymous reviewers for useful comments, and the students from the Fall 2008 Advanced Topics in Software Engineering class at our department for their feedback on this work. This material is based upon work partially supported by the NSF under Grant Nos. CCF-0746856, CNS-0615372, and CNS-0613665.

References

1. Beizer, B.: Software Testing Techniques. (1990)
2. Boyapati, C., Khurshid, S., Marinov, D.: Korat: automated testing based on Java predicates. In: ISSTA. (2002)
3. Sullivan, K., Yang, J., Coppit, D., Khurshid, S., Jackson, D.: Software assurance by bounded exhaustive testing. In: ISSTA. (2004)
4. Khurshid, S., Marinov, D.: TestEra: Specification-based testing of Java programs using SAT. *Auto. Soft. Eng. Jour.* (2004)
5. Daniel, B., Dig, D., Garcia, K., Marinov, D.: Automated testing of refactoring engines. In: ESEC/FSE. (2007)
6. Misailovic, S., Milicevic, A., Petrovic, N., Khurshid, S., Marinov, D.: Parallel test generation and execution with korat. In: ESEC/FSE. (2007)
7. Marinov, D., Andoni, A., Daniliuc, D., Khurshid, S., Rinard, M.: An evaluation of exhaustive testing for data structures. Technical report, MIT CSAIL (2003)

8. Jackson, D.: *Software Abstractions: Logic, Language and Analysis*. (2006)
9. Saff, D., Ernst, M.D.: Reducing wasted development time via continuous testing. In: *ISSRE*. (2003)
10. Do, H., Rothermel, G.: An empirical study of regression testing techniques incorporating context and lifetime factors and improved cost-benefit models. In: *ESEC/FSE*. (2006)
11. Stobie, K.: Model based testing in practice at Microsoft. *Electr. Notes Theor. Comput. Sci.* **111** (2005) 5–12
12. Rothermel, G., Untch, R.H., Chu, C., Harrold, M.J.: Test case prioritization: An empirical study. In: *ICSM*. (1999)
13. Elbaum, S., Malishevsky, A., Rothermel, G.: Incorporating varying test costs and fault severities into test case prioritization. In: *ICSE*. (2001)
14. Kim, J.M., Porter, A.: A history-based test prioritization technique for regression testing in resource constrained environments. In: *ICSE*. (2002)
15. Srivastava, A., Thiagarajan, J.: Effectively prioritizing tests in development environment. In: *ISSTA*. (2002)
16. Walcott, K.R., Soffa, M.L., Kapfhammer, G.M., Roos, R.S.: Time-aware test suite prioritization. In: *ISSTA*. (2006)
17. Yu, Y., Jones, J.A., Harrold, M.J.: An empirical study of the effects of test-suite reduction on fault localization. In: *ICSE*. (2008)
18. Rothermel, G., Elbaum, S., Malishevsky, A., Kallakuri, P., Davia, B.: The impact of test suite granularity on the cost-effectiveness of regression testing. In: *ICSE*. (2002)
19. Rothermel, G., Elbaum, S., Malishevsky, A.G., Kallakuri, P., Qiu, X.: On test suite composition and cost-effective regression testing. *ACM TOSEM* (2004)
20. Dickinson, W., Leon, D., Podgurski, A.: Finding failures by cluster analysis of execution profiles. In: *ICSE*. (2001)
21. Podgurski, A., Leon, D., Francis, P., Masri, W., Minch, M., Sun, J., Wang, B.: Automated support for classifying software failure reports. In: *ICSE*. (2003)
22. Liu, C., Yan, X., Fei, L., Han, J., Midkiff, S.P.: Sober: statistical model-based bug localization. In: *ESEC/FSE*. (2005)
23. Jones, J.A., Harrold, M.J., Bowring, J.F.: Debugging in parallel. In: *ISSTA*. (2007)
24. Runeson, P., Alexandersson, M., Nyholm, O.: Detection of duplicate defect reports using natural language processing. In: *ICSE*. (2007)
25. Liu, C., Zhang, X., Han, J., Zhang, Y., Bhargava, B.K.: Indexing noncrashing failures: A dynamic program slicing-based approach. In: *ICSM*. (2007)
26. Fowler, M., Beck, K., Brant, J., Opdyke, W., Roberts, D.: *Refactoring: Improving the Design of Existing Code*. (1999)
27. Eclipse Foundation, T.: Eclipse project. <http://www.eclipse.org>
28. Duncan, A.G., Hutchison, J.S.: Using attributed grammars to test designs and implementations. In: *ICSE*. (1981)
29. Maurer, P.M.: Generating test data with enhanced context-free grammars. *IEEE Soft.* (1990)
30. Zeller, A., Hildebrandt, R.: Simplifying and isolating failure-inducing input. *IEEE Trans. Soft. Eng.* (2002)
31. Rothermel, G., Harrold, M.J.: A safe, efficient regression test selection technique. *ACM TOSEM* (1997)
32. Heimdahl, M.P.E., Devaraj, G.: Test-suite reduction for model based tests: Effects on test quality and implications for testing. In: *ASE*. (2004)