

Directed Swarm Testing: Producing Random Regression Tests

Amin Alipour, Alex Groce, Rahul Gopinath, Arpit Christi

Abstract

While random testing can be a powerful and scalable method for finding faults in software, sophisticated random testers often test a whole program, not individual components. Writing such random testers for individual components of complex programs may require unreasonable effort. In this paper we present a novel method, *directed swarm testing*, that uses statistics and a variation of random testing to produce random tests that focus on only part of a program. One application is the ability to use an existing random tester to produce *regression tests* intended to detect faults in changed code. We demonstrate the effectiveness of this technique using real-world programs and test systems (the YAFFS2 file system, GCC, and Mozilla's SpiderMonkey JavaScript engine), and discuss various strategies for directed swarm testing. The best strategy can improve the coverage of targeted code by a factor ranging from 1.1-4.5x on average, and from nearly 3x to nearly 9x in the best case. For YAFFS2, directed swarm testing never decreased coverage, and for GCC and SpiderMonkey coverage increased for over 99% and 73% of targets, respectively, using the best strategy. Directed swarm testing improves detection of real SpiderMonkey faults. This lightweight technique is applicable to existing industrial-strength random testers.

1 Introduction

Random testing [18] (sometimes called fuzzing) is now widely recognized as an effective approach for testing software systems, including compilers [26, 30], standard libraries [25], static analysis systems [7], and file systems [15]. Random testing is used in both complex custom-built testing systems (such as those just cited) and simple test harnesses built in a couple of hours. Random testing is often easy to use, widely applicable, and can perform well in theory as well as practice [4]. However, random testing has a few important limitations. One critical limitation is that, for the most part, random testing is *ineffective* for regression testing. Random testers typically target an entire program or module, and have no mechanism for focusing testing on newly changed code (or any other kind of target of interest), other than writing a new test generator that focuses on the targeted code.

Much of the efficiency of random testing comes from its blind, undirected nature [31]. It is seldom practical to implement different random testers for all the potential focuses that might be needed, and the most powerful random testers [15, 26, 30]

tend to be based on generating complete inputs (e.g. programs or sequences of file system calls) as whole system tests for the Software Under Test (SUT), and do not even attempt to provide module-level testing. Of course, tests generated by a random tester can be used just like other tests in a regression suite, but replaying stored pre-existing tests defeats much of the point of random testing — the ability to produce novel inputs never considered by a developer and the ability to produce an essentially unlimited number of tests automatically, making effective use of any available testing budget and exploiting massive parallelism.

Techniques for making better use of random tests in regression testing are now appearing [13], but these do not allow the creation of true *random regression tests*: newly generated random tests that are specifically intended to test targeted (usually changed) code in a system. These are highly desirable for the simple reason that a large portion of newly modified code is buggy (perhaps up to one third [20]). Moreover, newly changed code has, by definition, been far less tested than long-standing code, especially in systems where aggressive random testing is applied routinely. The inability to perform efficient targeted testing is therefore a real deficiency in random testing.

While some other techniques (symbolic execution [11, 29] and search-based techniques [19, 23]) for test generation allow for targeting of specific source code, those techniques usually have not been scaled to the generation of, e.g., whole-program inputs for industrial strength compilers¹. Hand-tooled whole-program random testers, however, are a popular technique for testing such systems, including C compilers [22, 30], JavaScript engines [26], and Google’s Go language [27].

In this paper, we propose a method, *directed swarm testing*, that makes the generation of random tests targeting chosen code possible for many random testers, including all of the highly successful industrial-strength systems noted above. Using swarm testing [17], a variation of random testing, and recording statistical data on past testing results [16] enables generation of new random tests that target (that is, have higher probability of covering) any given source code element, usually without modifying an existing, highly-tuned random tester. This ability has further uses than just simple change-based regression testing; for example, a compiler developer using Csmith [30] and concerned about the correctness of a particular set of seldom-executed lines in a complex optimization’s implementation may apply this technique. Assuming that data on past testing has already been collected, the process can be as simple as putting the source lines of interest into a file and running a simple script that launches in parallel a large set of Csmith instances tuned to have high coverage of the suspect code. In our experiments, the fraction of tests that cover targeted code was improved by up to nearly 9x over running the random tester as usual, and the improvement is typically on the order of 2x or more. The more rarely code is covered in undirected tests — so long as it has been covered enough in past data to make a basis for statistical analysis — the more its coverage can be boosted.

Our experimental results show that, for single targets, across all strategies proposed, directed swarm testing improves the fraction of tests that hit a target by 3.5x on average

¹SAGE [10] has been applied to Internet Explorer’s Javascript engine, and performed better for code coverage (fault detection was not evaluated to our knowledge) than a mutation-based fuzzer and pure grammar-based random generation, but no comparison with a custom-tooled random tester such as jsfunfuzz has been published.

for YAFFS2, 2.5x on average for GCC, and 1.6x on average for SpiderMonkey. Directed swarm testing improved coverage for 100%, 95%, and 69.5% of targets (again, across all strategies) for YAFFS2, GCC, and SpiderMonkey respectively. Results for multiple targets are more complex, but still promising, though as the number of targets increases the effectiveness over swarm testing decreases (as it must, in the limit: targeting all code is equivalent to targeting none). We compare our method both against the baseline random test generators (hand-tooled optimized random testing) and the modified test generators using swarm testing.

Contributions of this paper include:

- A novel method (directed swarm testing) for generating random regression tests: new random tests that have greatly increased probability of covering selected source code targets (Section 3).
- Strategies for targeting both individual source code targets and multiple source code targets at once (Section 4).
- Empirical results showing the effectiveness of these strategies on large real-world software systems and test generators with complex test features (the YAFFS2 flash file system, the GCC compiler, and Mozilla’s SpiderMonkey JavaScript engine) (Sections 5 and 6).
- Empirical results of effectiveness of these strategies on finding *real faults* in a large software system (Sections 5 and 6).

2 Preliminary Concepts

2.1 Swarm Testing

Swarm testing [17] is a testing approach that improves the diversity of tests by randomizing the configuration of a test generation system (typically a random tester, though it is also applicable to model checking [1]). The idea behind swarm testing is simple: most significant random testing tools include some notion of *features*. A feature is a property of a test case that can be controlled by a test generator. A configuration of a test generator is often defined by a set of features. For example, in grammar-based testing, features are usually terminals or productions in the grammar, and in API-based testing each function or method call is a feature. The traditional approach to random testing is to always make all features available in the construction of each test, on the reasonable basis that if the configuration leaves out a feature, tests generated using that configuration cannot possibly detect faults in that feature’s behavior. Swarm testing, in contrast, randomly chooses (with base probability of 50%) which features to include in each test, omitting about half of all available features in each test. This often increases the effectiveness of testing due to interactions between features, and the fact that, since tests are limited in size, including many features necessarily means including less of each individual feature. For example, consider testing a stack implementation. Random tests that include both `push` and `pop` calls are extremely unlikely to detect overflow bugs that rely on too many `push` calls without a `pop`; tests omitting `pop` are almost

```

static uint16_t func_1(void) {
    uint16_t l_24[3][2] = {{0xD44FL, 0xD44FL},
        {0xD44FL, 0xD44FL}, {0xD44FL, 0xD44FL}};
    return l_24[1][1]; }
int main (int argc, char* argv[]) {
    func_1();
    return 0; }

```

(A) Simplified random test case for a C compiler, generated by Csmith, with boilerplate removed. Features here are C-language constructs. This test case features arrays but does not feature pointers, structs, jumps, volatiles, or 64 bit math.

```

tryItOut("L: {constructor = __parent__; }");
tryItOut("prototype = constructor;");
tryItOut("__proto__ = prototype;");
tryItOut("with({}){__proto__.__proto__=__parent__;}");

```

(B) Simplified random test case (without jsfunfuzz infrastructure) for SpiderMonkey JavaScript engine. Features here include labels, assignments, and with blocks, but do not include try blocks, infinite loops, or XML.

Figure 1: Features for Random Test Cases

guaranteed to find such a bug. Swarm testing has been recognized as essential to getting good results from compiler fuzzers [22] and in our own results has sometimes nearly *doubled* fault detection and/or coverage for mature random testers [16]. We (or our colleagues) have, beyond those results reported in previous work, also applied swarm testing with success to the CCG C compiler testing tool [24] and the GoSmith [27] fuzzer for Google’s Go language, with similar success. Le et al. report that using swarm testing was effective in their compiler validation efforts using equivalence modulo inputs [22].

Figure 1 shows examples of features for C and JavaScript-based test cases. Note that a feature can be a relatively simple grammatical construct or, depending on how tests are generated, a more complex semantic feature (e.g., irreducible control flow). Given a configuration, a hand-tuned random tester can usually generate a very large, or unbounded, number of different tests containing (at most) those features. For example, to use Csmith with a configuration is as easy as calling it with certain command line arguments (e.g., `csmith --no-pointers --no-structs --no-unions`).

2.2 Triggers and Suppressors

A *target* is any behavior of the SUT that is produced by some (but usually not all) test cases. The most obvious targets are faults and coverage entities, e.g.: whether a test case exposes a given fault, whether a given block or statement is executed, whether a branch is taken, or whether a particular path is followed. Hence, faults, blocks, branches, and paths are targets and a test case *hits* a target if it exposes or covers it. Given the concepts of features and targets, we can ask whether a feature f “helps” us to hit a target t : that is, are test cases with f more likely to hit t ? That some features are helpful for some targets is obvious: e.g., executing the first line of a method in an API library usually *requires* the call to be in the test case. Less obviously, features may make it *harder* to hit some targets. For example, finite-length tests of a bounded stack that contain `pop` calls are less likely to execute code that handles the case where the stack is full, closing files may make it harder to cover complex behavior in a file

system, and including pointers in a C program prevents some optimization passes from running [17].

There are three basic *roles* that a feature f can serve with respect to a target t : a *trigger*'s presence makes t easier to hit, a *suppressor*'s presence makes t harder to hit, and an irrelevant feature does not affect the probability of hitting t . The relation between features and targets can be non-trivial to predict and understand in large programs with complex features.

In previous work [16], we showed that for all non-trivial SUTs we examined, most targets had a few triggers and a few suppressors. We used a formal definition of trigger and suppressor based on Wilson scores [28] over hitting fractions in pure (undirected) swarm testing. Given feature f , target t , and test case population P where f appears in tests at rate r , compute a Wilson score interval for a given confidence (e.g., 95%) (l, h) on the true proportion of tests hitting t that contain f . If $h < r$, we can be, e.g., 95% confident that f suppresses t . The lower h is, the more suppressing f is for t . When $l > r$, f is a trigger for t . If neither of these cases holds, we can say that f is irrelevant to t . The appropriate bound (lower or upper) may then be used as a conservative estimate for the true fraction F of tests hitting t containing f :

$$F(f, t) = \begin{cases} r & \text{iff } l \leq r \leq h; & \text{(irrelevant)} \\ l & \text{iff } l > r; & \text{(trigger)} \\ h & \text{iff } h < r. & \text{(suppressor)} \end{cases}$$

F is easily interpreted when the rates for features are set at 50% in P , as in normal swarm testing. Critically, because of the way swarm testing works, feature/target relationships are always causal, evidence of a genuine semantic property of the SUT.

3 Directed Swarm Testing

We can exploit the fact that most targets of real-world SUTs have both triggers and suppressors to focus swarm testing on a given target, or set of targets. *Directed swarm testing* is performed similarly to conventional swarm testing, and like swarm testing, usually requires little or no modification of the base test generator. The difference between directed swarm testing and conventional swarm testing is that, instead of using completely random configurations, directed swarm testing uses configurations *based on the trigger and suppressor information collected for a single target or a set of targets*. Rather than a single algorithm, directed swarm testing is a family of strategies for choosing features in testing, with one constraint: when targeting t , directed swarm testing never uses configurations containing any suppressors of t .

When directed swarm testing is applied to multiple targets \mathcal{T} at once, as is often useful in testing changed code, it may only target some subset of \mathcal{T} in each individual test generation. A directed swarm testing strategy is effective if it increases the average rate at which tests hitting targets t are generated above the base rate for non-directed swarm testing. The larger the increase, the more effective the directed swarm testing strategy.

A typical application of directed swarm testing would be targeting changes made to the SUT. A developer has just implemented a new feature, and in the process added

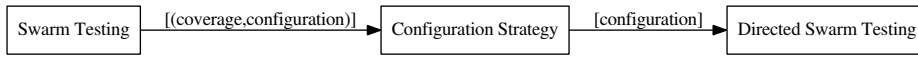


Figure 2: Workflow of directed swarm testing.

a new function f to the code, modified four lines of code in an existing function g , and added calls to f in three locations scattered throughout the program, all guarded by an existing conditional. The developer can run existing regression tests, and run an existing random tester in swarm mode, to detect bugs in the new feature. However, the function g is called by relatively few regression tests, and undirected swarm testing only calls g once in every twenty tests. The calls to f are only slightly more frequent. Assuming the unmodified code is correct, many of the tests generated in undirected swarm testing will be useless. However, it is easy to construct a set of targets for directed swarm testing in this situation: the four modified lines in g are obvious targets, and previous random testing results should contain enough information to calculate their triggers and suppressors with high confidence. The code for f , in contrast, is new; the developer has no information on triggers and suppressors for f itself. However, the developer always has information on some existing code that precedes new code to be targeted, and is as close as possible to it in the revised CFG for the SUT (the proof is trivial: if new code has no such nodes, it is either unreachable in the CFG, or the new code is the first node in the CFG, in which case it is always called and does not need to be targeted). The developer performs directed swarm testing, using this set of targets, and, if directed swarm testing is effective, is able to either find a bug or establish that the new code is likely correct much more quickly. The measure of success is how many tests covering changed code are produced within a given testing budget (or how quickly a fault is detected, when the code is faulty).

A major advantage of directed swarm testing is that, like swarm testing, it has essentially the same extremely low overhead as all random testing. The only additional cost for directed swarm testing is to collect coverage information when running some swarm tests, in order to compute triggers and successors for a program. Running some random tests with coverage instrumentation is already a common practice in aggressive testing, so this is hardly a major burden, even with the need to re-baseline trigger/suppressor information as code evolves over time. In practice, in our previous work, triggers and suppressors for lines of code that continued to exist through many software versions did not change dramatically, even from major release to major release, for Mozilla’s SpiderMonkey JavaScript engine [16]. Moreover, the cost of re-baselining is low enough that there is little reason not to routinely collect new trigger and suppressor information, especially as this is a by-product of an independently useful effort (checking code coverage).

4 Configuration Strategies

Figure 2 shows the overall workflow of directed swarm testing, which is simple. First, swarm testing is performed as usual, without any targets, to detect faults and collect

coverage information over the entire SUT. In order to apply directed swarm testing, the only information from this testing that is required is the set of (coverage, configuration) tuples for all tests generated in undirected swarm testing. This information can, as described in the introduction (Section 2.2) and in more detail in our previous work [16], be used to compute, for each source code target t (in this paper’s experiments, a statement), the set of triggering features $T(t)$, suppressing features $S(t)$, and irrelevant features $I(t)$. The heart of a directed swarm testing method is a strategy for producing configurations for new tests based on $T(t)$, $S(t)$, and $I(t)$. This can be done for a single t or for a set of targets \mathcal{T} . While the idea that knowledge of triggers and suppressors should enable us to improve testing for targets seems clear, there are trade-offs to consider in determining the actual configurations to use in testing for targets. Most importantly, the triggers and suppressors are determined with respect to a distribution of test cases such that most tests have about half of all features enabled; the causal patterns may not operate in the anticipated way when using a very different configuration distribution, due to the combinatoric complexities of hand-tuned random testers. While hitting the targets is important, it is also essential to maintain some test diversity to maximize the value of each individual test run — after all, simply running a single chosen test case that hits a target (with mutation fuzzing) may “maximize” target coverage, but loses almost all advantages of random testing.

4.1 Single-Target Strategies

We first consider the simplest case, targeting a single source code element. This is likely to be a very common goal, even for regression testing. If a developer only changes code in a single basic block, it is essentially one target with one set of triggers and suppressors (since the coverage vectors for all statements in a basic block are necessarily the same). Even modifying a few lines of code that are nearby in the CFG of the SUT is probably likely to involve similar triggers and suppressors, in most cases. In fact, multiple nearby targets can probably be effectively targeted in most cases by choosing their nearest common control flow dominator (for example, when all the modified code is in a single function).

We propose three basic strategies for a single target, t :

1. **Half-swarm:** The Half-swarm strategy produces configurations for testing in the same way as undirected swarm testing, with the exception that all features in $S(t)$ (the suppressors) are omitted from each configuration and all features in $T(t)$ are included in each configuration. It can be trivially implemented by applying an AND mask for suppressors (with all 1 bits except for suppressors) and an OR mask for triggers (with all 0 bits except for triggers) as a final stage in undirected swarm testing. In other words, a configuration $C_i = \{f | f \in S(t) \cup \text{random.Sample}(T(t))\}$, where *random.Sample* returns a random sample of a set such that each element has a 50% chance of being included.
2. **Triggers-only:** The Triggers-only strategy, as the name suggests, uses a single configuration for all testing, where all triggers are included and no other features are included: $C = \{f | f \in T(t)\}$.

3. **No-suppressors:** The No-suppressors strategy also uses only one configuration, which includes all triggers and irrelevant features, but no suppressors: $C = \{f \mid f \notin S(t)\}$.

The motivation for **Half-swarm** is that swarm testing is effective, and directed swarm testing should, perhaps, remain as close to undirected swarm testing as possible, except for taking triggers and suppressors into account. The motivation for the other two strategies is that while swarm testing is effective for general testing of an SUT, it may not be ideal when random testing of a single target is the goal. The diversity that makes swarm testing useful may be useless or harmful with a single target; however, it is not clear if a minimal or maximal configuration that respects triggers and suppressors would be best, given this assumption. **Triggers-only** uses a minimal configuration, with only those features known to improve coverage of the target included, while **No-suppressors** is maximal, only omitting features known to hinder coverage of the target. The computational cost for all techniques is the same (and essentially identical to that of non-directed swarm testing or pure random testing). As we see below, in addition to the basic empirical question of effectiveness, the idiosyncracies of some random testers may also determine which of these strategies should be chosen. In particular, for some testers, if very few features are present in a configuration, it may not generate any valid tests. When there are many features and a 50% chance of inclusion, the problem does not arise, but using Triggers-only may frequently fail.

4.2 Multiple-Target Strategies

For multiple targets, \mathcal{T} , our strategies reduce the problem to that for single targets $t \in \mathcal{T}$:

1. **Round-robin:** The Round-robin strategy simply applies a single-target strategy in a round-robin fashion, for $t \in \mathcal{T}$.
2. **Merging:** The Merging strategy attempts to *merge* triggers and suppressors for targets in \mathcal{T} to produce a minimal set of targets (each of which may represent multiple real targets) then uses round-robin.

The motivation behind **Round-robin** is simple: to cover a set of targets, split the testing time between those targets. If multiple targets have similar suppressors and triggers, we may end up covering a target with tests not aimed at that target, but the basic idea is simply to assume all targets are equally important and cannot be tested at once. Round-robin is parameterized on a single-target strategy.

Merging approaches are more complex. They are motivated by an observation: if for two targets, t_1 and t_2 , $\neg \exists f. (f \in S(t_1) \wedge f \in T(t_2)) \vee (f \in T(t_1) \wedge f \in S(t_2))$, then there may be no reason we have to target t_1 and t_2 with different configurations. They do not have any *conflicts*, where a conflict is a feature that suppresses one target but triggers the other target.

Algorithm 1 illustrates one simple algorithm to produce a set of targets \mathcal{T}' for targets \mathcal{T} . Given targets $t_i, t_j \in \mathcal{T}$, we say t_j *subsumes* t_i , denoted $t_j \sqsupseteq t_i$, if and only if, $S(t_i) \subset S(t_j) \wedge T(t_i) \subset T(t_j)$. In other words, t_j requires a *stricter* combination of

features than t_i . **Subsumption** merging removes t_i and only keeps the stricter combination of features, assuming that it will test both targets. The computational cost of the algorithm is quadratic in the number of targets to consider merging (and thus negligible for likely sets of targets).

Algorithm 1 Algorithm for Merging using Subsumption only.

```

1: for  $\forall t_i \in \mathcal{T}$  do
2:   if  $\exists t_j \in \mathcal{T} | t_j \sqsupset t_i$  then
3:      $\mathcal{D} = \mathcal{D} \cup t_i$ 
4:   end if
5: end for
6: return  $t \in \mathcal{T} | t \notin \mathcal{D}$ 

```

Algorithm 2 Algorithm for Aggressive Merging, with randomized approximation of optimal merges ($n = \#$ of trials).

```

1:  $\mathcal{B} = \mathcal{T}$ 
2: for  $i = 0 \dots n - 1$  do
3:    $\mathcal{M} = \mathcal{T}$ 
4:   while  $\exists t_i, t_j \in \mathcal{M} : t_i \neq t_j \wedge (\neg \exists f. (f \in S(t_i) \wedge f \in T(t_j)) \vee (f \in T(t_i) \wedge f \in S(t_j)))$  do
5:     pick  $t_i, t_j$ 
6:      $T(t_m) = f | f \in T(t_i) \vee f \in T(t_j)$ 
7:      $S(t_m) = f | f \in S(t_i) \vee f \in S(t_j)$ 
8:      $I(t_m) = f | f \notin T(t_m) \wedge f \notin S(t_j)$ 
9:      $\mathcal{M} = \mathcal{M} \cup t_m - t_i - t_j$ 
10:  end while
11:  if  $|\mathcal{M}| < |\mathcal{B}|$  then
12:     $\mathcal{B} = \mathcal{M}$ 
13:  end if
14: end for
15: return  $\mathcal{B}$ 

```

It is also possible to merge in a more **Aggressive** fashion. In the absence of conflicts, we can in principle merge *any* two targets even where neither is stricter than the other, treating them as one target t' , with $T(t') = f | f \in T(t_1) \vee f \in T(t_2)$, $S(t') = f | f \in S(t_1) \vee f \in S(t_2)$, and $I(t') = f | f \notin S(t') \wedge f \notin T(t')$. In this way, we can keep merging targets (replacing the two non-conflicting targets with the new meta-target) without conflicts to produce a small set of configurations that are directed at many targets at once. However, finding the merges to produce a truly minimal set of configurations is, we believe, an NP-complete problem, since non-subsuming merges are path-dependent (e.g., c_i may be merged with either c_j or c_k but not both, affecting future merges). We implemented an SMT-based exact solver for merging targets using Z3 [8], which was able to construct perfect solutions for up to 20 targets (typically solving for 300 features in less than 2 minutes, but sometimes taking more than 10

minutes), but did not scale to 40 targets at all, even with very few features (timing out after many hours). Fortunately, due to the fact that most targets have either absolutely few (< 3) triggers and suppressors or at least relatively few ($< 5\%$ of features) triggers and suppressors [16], random ordering of matches (using the best solution after a fixed number of trials) approximates exact solutions effectively and quickly. In our experiments, a random approximation of optimal merging, even using 1,000 trials, always produced a nearly-optimal set of configurations (at most one larger than the optimal set produced by Z3) in less than 1 second, for up to 20 targets. In experiments, we used 10,000 trials. Algorithm 2 shows the randomized algorithm for Aggressive Merging of targets. We assume that Subsumption Merging has already been applied before this algorithm is called.

Both the Subsumption and Aggressive Merging strategies are, like the Round-robin strategy, parameterized on a single-target configuration strategy. It is, in part for this reason, not clear whether (and how much) we *should* merge configurations. Merging targets produces “more specialized” configurations that leave little room for the basic single-target strategies to operate (because merging increases the numbers of fixed triggers and suppressors for each merged target). Round-robin maintains maximal configuration diversity (consistent with directing the testing). Subsumption Merging assumes that when one target subsumes another, they are truly similar and can be tested in the same way. Aggressive Merging uses as few configurations as possible, but may result in a very small number of targets with very few irrelevant features. Whether such targets can actually be effectively tested by the same configurations is not obvious without empirical investigation.

5 Evaluation Methodology

We used three medium-moderately large C programs (shown in Table 1) to evaluate directed swarm testing. While not extremely large, these are all important systems-software programs, the typical of the kind of program for which an effective dedicated random tester can be expected to exist. For YAFFS2 (formerly the default image file system for Android), we used our own in-house test generation tools, which are descended from those used to test the file systems for NASA’s Curiosity Mars Rover [15]. For GCC, we used the Csmith [30] C compiler fuzzer to generate tests. Csmith is a highly effective tool that has been used to detect more than 400 previously unknown bugs in GCC, LLVM, and other production C compilers. For SpiderMonkey, Mozilla’s JavaScript engine, we used `jsfunfuzz` [26], a well-known JavaScript fuzzer responsible for finding more than 6,400 bugs in SpiderMonkey. We had previously modified `jsfunfuzz` to support swarm testing. The other two test generators already supported swarm testing.

Table 2 shows parameters for our experiments. In this table: # Features shows the number of features in the SUT that can be tested by the corresponding fuzzer, seed time shows time spent in minutes to generate the initial (undirected swarm) test suite that is used for extracting trigger/suppressor features for statements, and directed time shows the time spent for directed testing of targets. The stochastic nature of random testing required us to run experiments multiple times to ensure results are statistically

Table 1: Experimental Subjects

SUT	LOC	Fuzzer	Description
YAFFS2	15K	yaffs2tester	Flash File System
GCC 4.4.7	860K	Csmith	C and C++ Compiler
SpiderMonkey 1.6	118K	jsfunfuzz	JavaScript Engine For Mozilla

Table 2: Experimental Parameters.

SUT	# Features	Seed time (min.)	Directed time (min.)	# Undirected Suites
YAFFS2	43	15	5	60
GCC	25	60	10	30
SpiderMonkey	171	30	10	54

significant. For each test subject we generated between 30 and 60 initial test suites (# Suites) using undirected swarm testing. We collected data on configurations and coverage from these tests, and computed Wilson scores (and thus triggers and suppressors) for all statements covered in the tests. For each such test suite, we picked up to 35 sets of random targets (statements), with sizes 1, 5, 10 and 20 (up to 5 for each size) to evaluate directed swarm testing². Moreover, we used the default configuration of each test generator to generate one traditional (non-swarm) random test suite for each swarm test suite, to further compare the effectiveness of directed swarm testing and traditional random testing.

We randomly chose targets (statements) covered by 10% to 30% of test cases in the original test suite, to restrict evaluation to targets that are at least somewhat difficult to cover, but for which a statistical basis for directed swarm testing definitely exists. For more rarely covered targets, where triggers and suppressors are less certain, the nearest control-flow dominator with sufficient coverage in tests can be used as a replacement target. Note that with a large amount of historical coverage data, as might be collected in an overnight test run on a stable version, many more targets would have statistical support for accurate triggers and suppressors. The 10%-30% selection is only to enable experiments using limited coverage data, not a limitation of directed swarm testing.

For the single-target sets we applied each of the Half-swarm, Triggers-only, and No-suppressors strategies. For all multiple target sets, we also applied Round-robin, Subsumption, and Aggressive strategies (in each case paired with a single-target strategy, for nine strategies in all). We varied the time for undirected testing and directed testing according to suite complexity in each case. In total, we ran tests for slightly more than 3,000 hours and generated over 20,00 directed test suites.

²We also collected data for size 2, 3, and 4 regressions, which will be provided in a technical report; in the interests of space, these results, which shed little light on multi-target strategies and were similar to results for size 5, are omitted from this version.

Our primary measure of effectiveness is simple. For any test suite, we compute the hitting fraction HF for tests that cover a target t (if there are n tests in a suite and m tests cover t then, $HF = \frac{m}{n}$) — if every test in a suite covers t , $HF = 1.0$ and if no tests cover t , $HF = 0.0$. Suppose the hitting ratios in undirected test and directed test are HF_u and HF_d respectively, we use $\frac{HF_d}{HF_u}$ ratio to measure the effectiveness of directed testing in hitting targets. Note that directed test suites with $\frac{HF_d}{HF_u} > 1.0$ offer improvement over the undirected test suites. This is the measure a developer wants to increase when targeting a particular statement.

6 Results

Our experimental results address six basic research questions:

- **RQ1:** (How much) does directed swarm testing improve coverage for single targets?
- **RQ2:** Which strategies for single-target directed swarm testing are most effective?
- **RQ3:** (How much) does directed swarm testing improve coverage for multiple targets at once?
- **RQ4:** Which strategies for multiple-target directed swarm testing are most effective?
- **RQ5:** Does directed swarm testing help detect actual faults?
- **RQ6:** How much does directed swarm testing improve coverage over traditional random testing?

Figure 3 illustrates the distribution of targets’ hitting fraction (HF) for (undirected) swarm testing and directed swarm testing. It shows that, in most cases, the hitting fraction for targets in directed swarm testing is much higher than the hitting fraction for undirected swarm testing. For brevity, in the rest of this section, we use “directed swarm testing” and “directed testing” interchangeably, as directed swarm testing is the only directed testing approach we evaluate (and, to our knowledge, the only one applicable to our subject programs).

Tables 3-6 provide much more detailed information about the performance of directed testing. These summarize the $\frac{HF_d}{HF_u}$ under different strategies. In these tables, In this table, “count” row contains the number of test suites generated by directed swarm testing using strategies described in the corresponding row. $\frac{HF_d}{HF_u} > 1.0$ column shows the fraction of test suites (i.e. count) where target(s) were covered more often by the directed swarm testing than the corresponding initial undirected swarm testing. For example, the value 0.8 in this column means: in 80% of test suites generated by directed swarm testing, the HF for targets is higher than the original undirected swarm. “mean”, “std. dev”, “min”, “25%”, “50%”, “75%” and, “max” respectively denote average, standard deviation, minimum, first quartile, second quartile (i.e. median), third quartile and maximum of $\frac{HF_d}{HF_u}$ in test suites generated by corresponding strategies in each row.

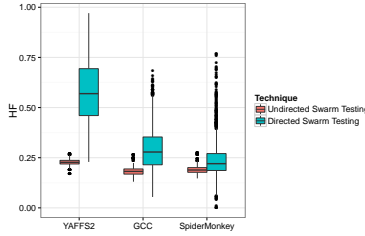


Figure 3: Hitting fraction in undirected swarm testing (HF_u) versus directed swarm testing (HF_d) over all strategies.

Table 3: Results for single-target directed random testing.

Strategy	$\frac{HF_d}{HF_u} > 1$	count	mean	std. dev	min	25%	50%	75%	max	p-val
YAFFS2										
Half-swarm	1.0	218.0	3.56	0.59	1.38	3.11	3.7	3.96	5.01	0.0
Triggers-only	1.0	218.0	3.94	0.64	2.26	3.57	4.0	4.25	7.87	0.0
No-suppressors	1.0	216.0	3.03	0.7	1.03	2.4	3.19	3.56	4.44	0.0
GCC										
Half-swarm	0.99	138.0	2.4	0.99	0.94	1.69	2.19	3.0	6.33	0.0
Triggers-only	0.92	129.0	2.28	1.0	0.53	1.53	2.13	2.94	5.29	0.0
No-suppressors	0.94	135.0	2.56	1.0	0.0	1.86	2.59	3.23	5.58	0.0
SpiderMonkey										
Half-swarm	0.73	260.0	1.75	1.06	0.0	0.88	1.74	2.49	4.39	0.0
Triggers-only	0.84	19.0	4.56	3.01	0.11	2.6	3.62	7.23	8.82	0.0006
No-suppressors	0.65	260.0	1.15	0.62	0.0	0.61	1.27	1.6	3.14	0.30234

6.1 RQ1 and RQ2: Single-Target Strategies

Table 3 shows the results for single-target directed swarm testing under different directed testing strategies, including p -values for Wilcoxon tests. Figure 4 visualizes these results. Table 3 shows that directed swarm testing has been successful in increasing HF of targets in YAFFS and GCC. For all cases in YAFFS, directed swarm testing increased the hitting ratio of the targets. The hitting fraction of targets using directed swarm testing was more than three times more than the hitting fraction in the undirected testing, on average. For GCC, directed swarm testing could increase the hitting fraction of targets for more than 90% of targets. On average, the directed testing increased the hitting fraction of targets by a factor of 2 or more.

The results for SpiderMonkey are mixed partly because the design of `jsfunfuzz` is such that, if we remove certain features, the fuzzer cannot produce any test cases at all. Moreover `jsfunfuzz` encodes SpiderMonkey’s feature set by paths through a complex recursive code generation system that resembles a grammar. In many cases, with SpiderMonkey, the triggers for a target are low-level productions that are only reachable through top-level parts of the fuzzer that correspond to irrelevant features — they are highly redundant. This makes it hard to identify triggers and suppressors, since the chance of undirected swarm generating a configuration disabling all paths is small. However, even for SpiderMonkey, directed swarm testing increases the hitting

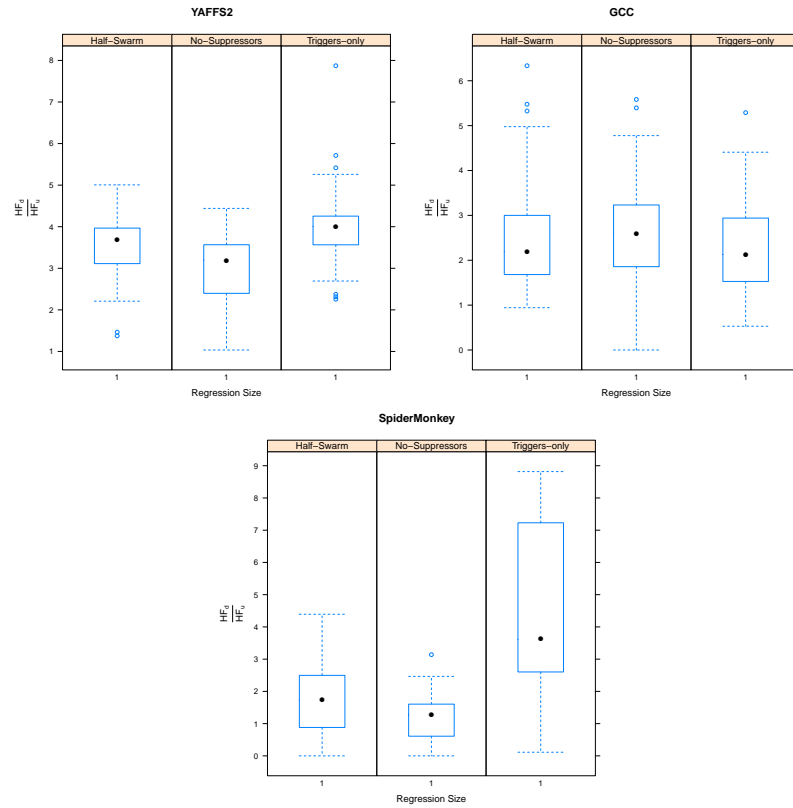


Figure 4: Single-target strategies compared.

fraction of more than half of targets, and Half-swarm had mean improvement of close to 2x. Note that most configurations for the Triggers-only strategy could not generate any test cases.

Observation 1: Directed swarm testing, with the exception of one strategy for SpiderMonkey, significantly ($p < 0.01$) increases coverage over undirected testing.

The average improvement for single-target directed testing ranges from 1.15x for SpiderMonkey with the No-Suppressor strategy to nearly 4x for Triggers-only with YAFFS2.

Observation 2: There is no clear best strategy for single-target testing, though it is clear that adopting Triggers-only may be risky in some settings.

6.2 RQ3 and RQ4: Multiple-Target Strategies

For analysis of multi-target directed testing, we use the *average* hitting fraction, i.e. \overline{HF} , for comparison of effectiveness of directed testing. Figure 5 shows results $\frac{HF_d}{HF_u}$ with various regression sizes. The most obvious trend is that, while it is effective, the effectiveness of directed testing decreases with an increase in number of targets. For YAFFS, targeted testing always increases the hitting fraction of targets, on average between 1.89x to 3.77x .

In GCC, directed testing improves hitting fraction for 87.2% of all target sets. No-suppressors and Half-swarm strategies improve the hitting fraction of targets in 98.2% of cases. On average, they improve hitting factor between 1.35x to 2.87x.

Directed swarm testing improves the hitting fraction for 74% of target sets in SpiderMonkey. The Triggers-only strategy for SpiderMonkey could not generate tests for many targets, due to the complex recursive code generation in `jsfunfuzz` (mentioned earlier) generating test suites for only about 30% of targets. Half-swarm and No-suppressor strategies improve the hitting fraction 73.9% of their targets, between 1.06x and 1.69x on average, for regression sizes of 4 and 2, respectively.

Tables 4–6 show detailed results for multiple-target testing. The second lesson here is that for all subjects, Triggers-only was ineffective. Given the risks seen in single-target testing and the lackluster results here, we believe that Triggers-only may be the least effective strategy, despite its good results for YAFFS2 single-target directed swarm testing. It may be that Triggers-only is simply too extreme: conventional random testing uses all features in every test, and swarm testing can often improve this by reducing the fraction of features by half. Lowering it to the small number of triggers for many targets may simply not match the behavior most random testers are designed to work with, or produce too little complex interaction of software components to provide good testing.

Observation 3: For YAFFS, GCC, and SpiderMonkey, for No-suppressor and Half-swarm strategies, the hitting fraction of at least 95% of target sets increases significantly using directed swarm testing ($p < 0.01$, Wilcoxon rank-sum).

Figure 6 shows the performance of different merge strategies across test subjects, and Figure 7 shows how merge strategies affected the number of effective targets (how much merging was possible). Aggressive merging produced consistently very small sets of targets, while Subsumption results are generally closer to Round-robin than to Aggressive. The difference between the Round-robin and Subsumption strategies in hitting fractions was therefore minimal (and in most cases, not statistically significant). The most likely explanation is that when two targets have similar enough triggers and suppressors to merge, for our subjects, testing one target in round robin is likely to “accidentally” target the other target as well. Aggressive merging improved hitting fractions for GCC and SpiderMonkey, but performed poorly for YAFFS2.

6.3 RQ5: Actual Fault Detection

In addition to our basic results showing that directed swarm testing can improve the *coverage* of targets, we also performed experiments on actual fault detection. These

Table 4: Detailed multi-target results.

Regression Size	Merge Mode	Strategy	$\frac{HFE_c}{HFE_n} > 1.0$	count	mean	std. dev	min	25%	50%	75%	max	p-val
YAFFS2												
2	Round-robin	Half-swarm	1.0	52.0	3.22	0.59	1.83	2.98	3.32	3.65	4.21	0.0
2	Round-robin	Triggers-only	1.0	52.0	3.42	0.79	1.65	2.99	3.61	3.97	4.97	0.0
2	Round-robin	No-suppressors	1.0	50.0	2.93	0.46	1.49	2.68	2.94	3.24	3.75	0.0
2	Subsumption	Half-swarm	1.0	88.0	3.18	0.68	1.78	2.64	3.33	3.69	5.31	0.0
2	Subsumption	Triggers-only	1.0	88.0	3.38	0.96	1.57	2.64	3.59	3.96	7.86	0.0
2	Subsumption	No-suppressors	1.0	88.0	2.91	0.51	1.55	2.59	2.95	3.28	3.88	0.0
2	Aggressive	Half-swarm	1.0	53.0	3.44	0.43	2.04	3.21	3.53	3.74	4.16	0.0
2	Aggressive	Triggers-only	1.0	53.0	3.77	0.48	2.24	3.47	3.79	4.03	5.14	0.0
2	Aggressive	No-suppressors	1.0	52.0	2.97	0.48	1.57	2.79	2.98	3.26	3.91	0.0
3	Round-robin	Half-swarm	1.0	65.0	2.87	0.65	1.55	2.35	2.99	3.43	3.9	0.0
3	Round-robin	Triggers-only	1.0	65.0	2.95	0.87	1.31	2.2	3.05	3.82	4.22	0.0
3	Round-robin	No-suppressors	1.0	65.0	2.83	0.52	1.25	2.61	2.93	3.21	3.78	0.0
3	Subsumption	Half-swarm	1.0	76.0	2.75	0.69	1.33	2.17	2.58	3.38	3.9	0.0
3	Subsumption	Triggers-only	1.0	76.0	2.79	0.91	1.3	2.11	2.43	3.72	4.25	0.0
3	Subsumption	No-suppressors	1.0	76.0	2.75	0.57	1.26	2.49	2.87	3.16	3.76	0.0
3	Aggressive	Half-swarm	1.0	65.0	2.98	0.71	1.5	2.28	3.32	3.58	3.96	0.0
3	Aggressive	Triggers-only	1.0	65.0	3.15	0.89	1.52	2.11	3.59	3.9	4.49	0.0
3	Aggressive	No-suppressors	1.0	65.0	2.84	0.6	1.18	2.51	3.0	3.25	3.79	0.0
4	Round-robin	Half-swarm	1.0	67.0	2.79	0.58	1.56	2.4	2.69	3.26	3.98	0.0
4	Round-robin	Triggers-only	1.0	67.0	2.86	0.78	1.29	2.29	2.65	3.56	4.21	0.0
4	Round-robin	No-suppressors	1.0	67.0	2.77	0.46	1.53	2.45	2.79	3.09	3.79	0.0
4	Subsumption	Half-swarm	1.0	71.0	2.71	0.62	1.45	2.28	2.59	3.26	3.92	0.0
4	Subsumption	Triggers-only	1.0	71.0	2.75	0.82	1.31	2.06	2.53	3.53	4.14	0.0
4	Subsumption	No-suppressors	1.0	71.0	2.72	0.49	1.55	2.44	2.73	3.08	3.66	0.0
4	Aggressive	Half-swarm	1.0	67.0	2.9	0.69	1.57	2.24	3.13	3.49	3.99	0.0
4	Aggressive	Triggers-only	1.0	67.0	3.03	0.85	1.32	2.14	3.35	3.75	4.14	0.0
4	Aggressive	No-suppressors	1.0	67.0	2.77	0.57	1.41	2.56	2.89	3.17	3.82	0.0
5	Round-robin	Half-swarm	1.0	69.0	2.69	0.57	1.48	2.23	2.61	3.13	4.02	0.0
5	Round-robin	Triggers-only	0.99	69.0	2.65	0.73	1.0	2.1	2.55	3.15	4.51	0.0
5	Round-robin	No-suppressors	1.0	69.0	2.82	0.45	1.55	2.54	2.87	3.11	3.61	0.0
5	Subsumption	Half-swarm	1.0	72.0	2.65	0.6	1.51	2.17	2.53	3.15	4.07	0.0
5	Subsumption	Triggers-only	0.99	72.0	2.61	0.77	0.99	2.07	2.42	3.29	4.48	0.0
5	Subsumption	No-suppressors	1.0	72.0	2.79	0.47	1.43	2.49	2.81	3.09	3.67	0.0
5	Aggressive	Half-swarm	1.0	70.0	2.68	0.78	1.54	1.99	2.34	3.48	4.28	0.0
5	Aggressive	Triggers-only	1.0	69.0	2.7	0.92	1.34	1.97	2.15	3.63	4.48	0.0
5	Aggressive	No-suppressors	1.0	69.0	2.68	0.66	1.28	2.24	2.81	3.2	3.86	0.0
10	Round-robin	Half-swarm	1.0	71.0	2.53	0.48	1.5	2.16	2.54	2.85	3.56	0.0
10	Round-robin	Triggers-only	1.0	71.0	2.47	0.62	1.28	2.0	2.46	2.82	3.88	0.0
10	Round-robin	No-suppressors	1.0	71.0	2.66	0.38	1.3	2.39	2.72	2.88	3.38	0.0
10	Subsumption	Half-swarm	1.0	71.0	2.42	0.5	1.53	2.03	2.41	2.69	3.59	0.0
10	Subsumption	Triggers-only	1.0	71.0	2.33	0.65	1.16	1.82	2.37	2.62	3.84	0.0
10	Subsumption	No-suppressors	1.0	71.0	2.64	0.37	1.74	2.34	2.63	2.91	3.41	0.0
10	Aggressive	Half-swarm	1.0	71.0	2.2	0.65	1.33	1.78	1.99	2.45	3.72	0.0
10	Aggressive	Triggers-only	1.0	71.0	2.16	0.76	1.27	1.61	1.98	2.36	4.0	0.0
10	Aggressive	No-suppressors	1.0	71.0	2.35	0.62	1.27	1.79	2.25	2.85	3.53	0.0
20	Round-robin	Half-swarm	1.0	74.0	2.53	0.33	1.93	2.25	2.5	2.75	3.36	0.0
20	Round-robin	Triggers-only	1.0	74.0	2.47	0.47	1.53	2.12	2.38	2.86	3.7	0.0
20	Round-robin	No-suppressors	1.0	74.0	2.69	0.2	2.24	2.57	2.69	2.82	3.14	0.0
20	Subsumption	Half-swarm	1.0	76.0	2.33	0.32	1.8	2.1	2.27	2.53	3.35	0.0
20	Subsumption	Triggers-only	1.0	76.0	2.22	0.45	1.47	1.87	2.13	2.48	3.68	0.0
20	Subsumption	No-suppressors	1.0	76.0	2.6	0.3	1.03	2.45	2.61	2.77	3.14	0.0
20	Aggressive	Half-swarm	1.0	75.0	1.97	0.49	1.23	1.69	1.9	2.17	3.63	0.0
20	Aggressive	Triggers-only	1.0	75.0	1.89	0.57	1.19	1.52	1.78	2.07	3.67	0.0
20	Aggressive	No-suppressors	1.0	74.0	2.18	0.52	1.21	1.82	2.09	2.59	3.55	0.0

Table 5: Detailed multi-target results.

Regression Size	Merge Mode	Strategy	$\frac{HFE}{HFE_{\text{ref}}} > 1.0$	count	mean	std. dev	min	25%	50%	75%	max	p-val
GCC												
2	Round-robin	Half-swarm	1.0	4.0	2.05	0.19	1.81	1.95	2.07	2.17	2.25	0.03394
2	Round-robin	Triggers-only	1.0	4.0	1.66	0.57	1.03	1.28	1.65	2.02	2.31	0.03394
2	Round-robin	No-suppressors	1.0	4.0	2.57	0.25	2.21	2.5	2.67	2.74	2.75	0.03394
2	Subsumption	Half-swarm	1.0	6.0	1.89	0.3	1.45	1.73	1.95	2.02	2.31	0.01385
2	Subsumption	Triggers-only	0.67	6.0	1.24	0.5	0.4	1.06	1.38	1.5	1.82	0.12443
2	Subsumption	No-suppressors	1.0	6.0	2.54	0.49	1.99	2.18	2.47	2.81	3.3	0.01385
2	Aggressive	Half-swarm	1.0	4.0	2.65	0.73	1.89	2.33	2.53	2.84	3.64	0.03394
2	Aggressive	Triggers-only	1.0	4.0	2.06	0.39	1.67	1.74	2.08	2.39	2.4	0.03394
2	Aggressive	No-suppressors	1.0	4.0	2.87	0.42	2.33	2.64	2.93	3.15	3.28	0.03394
3	Round-robin	Half-swarm	1.0	4.0	1.45	0.49	1.08	1.11	1.28	1.62	2.13	0.03394
3	Round-robin	Triggers-only	0.75	4.0	1.0	0.07	0.89	0.99	1.03	1.04	1.06	0.3575
3	Round-robin	No-suppressors	1.0	4.0	2.06	0.44	1.59	1.87	2.01	2.19	2.64	0.03394
3	Subsumption	Half-swarm	1.0	4.0	1.5	0.41	1.11	1.23	1.44	1.71	2.03	0.03394
3	Subsumption	Triggers-only	0.5	4.0	1.19	0.54	0.7	0.85	1.07	1.41	1.93	0.6425
3	Subsumption	No-suppressors	1.0	4.0	2.08	0.35	1.66	1.86	2.13	2.35	2.41	0.03394
3	Aggressive	Half-swarm	1.0	4.0	2.17	0.67	1.58	1.59	2.16	2.74	2.79	0.03394
3	Aggressive	Triggers-only	1.0	4.0	1.61	0.61	1.04	1.12	1.58	2.07	2.23	0.03394
3	Aggressive	No-suppressors	1.0	4.0	2.33	0.86	1.48	1.68	2.28	2.93	3.28	0.03394
4	Round-robin	Half-swarm	0.8	5.0	1.59	0.51	0.85	1.41	1.68	1.75	2.26	0.03981
4	Round-robin	Triggers-only	0.4	5.0	1.14	0.55	0.7	0.89	0.97	1.01	2.1	0.65708
4	Round-robin	No-suppressors	1.0	5.0	2.36	0.59	1.68	1.9	2.48	2.59	3.16	0.02156
4	Subsumption	Half-swarm	1.0	5.0	1.74	0.4	1.44	1.49	1.59	1.78	2.42	0.02156
4	Subsumption	Triggers-only	0.4	5.0	1.03	0.39	0.58	0.9	0.92	1.13	1.63	0.55363
4	Subsumption	No-suppressors	1.0	5.0	2.39	0.79	1.51	1.88	2.39	2.6	3.59	0.02156
4	Aggressive	Half-swarm	1.0	5.0	1.96	0.44	1.25	1.99	2.01	2.12	2.44	0.02156
4	Aggressive	Triggers-only	1.0	5.0	2.06	0.3	1.59	2.02	2.08	2.23	2.39	0.02156
4	Aggressive	No-suppressors	1.0	5.0	2.69	0.58	1.98	2.31	2.74	2.97	3.47	0.02156
5	Round-robin	Half-swarm	1.0	4.0	1.41	0.13	1.3	1.3	1.39	1.49	1.55	0.03394
5	Round-robin	Triggers-only	0.75	4.0	1.17	0.16	0.95	1.1	1.21	1.28	1.33	0.07206
5	Round-robin	No-suppressors	1.0	4.0	2.02	0.23	1.82	1.9	1.96	2.09	2.34	0.03394
5	Subsumption	Half-swarm	1.0	4.0	1.43	0.2	1.21	1.3	1.4	1.53	1.68	0.03394
5	Subsumption	Triggers-only	0.75	4.0	1.26	0.34	0.79	1.16	1.33	1.44	1.59	0.07206
5	Subsumption	No-suppressors	1.0	4.0	1.88	0.44	1.25	1.77	2.04	2.16	2.19	0.03394
5	Aggressive	Half-swarm	1.0	4.0	1.46	0.2	1.26	1.3	1.47	1.62	1.64	0.03394
5	Aggressive	Triggers-only	1.0	4.0	1.56	0.21	1.26	1.53	1.64	1.67	1.71	0.03394
5	Aggressive	No-suppressors	1.0	4.0	1.9	0.17	1.78	1.78	1.84	1.96	2.14	0.03394
10	Round-robin	Half-swarm	1.0	5.0	1.45	0.29	1.12	1.35	1.42	1.44	1.91	0.02156
10	Round-robin	Triggers-only	0.4	5.0	0.85	0.33	0.56	0.61	0.73	1.04	1.33	0.82738
10	Round-robin	No-suppressors	1.0	5.0	2.12	0.21	1.92	1.95	2.05	2.28	2.38	0.02156
10	Subsumption	Half-swarm	0.83	6.0	1.38	0.28	0.92	1.34	1.39	1.44	1.81	0.0232
10	Subsumption	Triggers-only	0.33	6.0	0.93	0.26	0.43	0.96	0.98	1.0	1.23	0.82728
10	Subsumption	No-suppressors	1.0	5.0	2.08	0.31	1.65	1.93	2.07	2.34	2.43	0.02156
10	Aggressive	Half-swarm	1.0	5.0	2.1	0.57	1.45	1.93	2.07	2.07	3.01	0.02156
10	Aggressive	Triggers-only	1.0	5.0	1.66	0.37	1.32	1.5	1.57	1.62	2.28	0.02156
10	Aggressive	No-suppressors	1.0	5.0	2.28	0.41	1.68	2.15	2.25	2.58	2.75	0.02156
20	Round-robin	Half-swarm	0.8	5.0	1.44	0.34	0.92	1.29	1.55	1.71	1.72	0.03981
20	Round-robin	Triggers-only	0.2	5.0	0.85	0.18	0.68	0.73	0.82	0.87	1.15	0.93099
20	Round-robin	No-suppressors	1.0	5.0	1.91	0.29	1.57	1.73	1.84	2.21	2.21	0.02156
20	Subsumption	Half-swarm	1.0	5.0	1.35	0.22	1.06	1.22	1.35	1.49	1.62	0.02156
20	Subsumption	Triggers-only	0.0	5.0	0.78	0.13	0.62	0.73	0.78	0.81	0.98	0.97844
20	Subsumption	No-suppressors	1.0	5.0	1.93	0.18	1.69	1.88	1.91	1.98	2.18	0.02156
20	Aggressive	Half-swarm	1.0	5.0	1.73	0.34	1.4	1.47	1.67	1.87	2.24	0.02156
20	Aggressive	Triggers-only	1.0	5.0	1.66	0.29	1.31	1.48	1.64	1.95	1.95	0.02156
20	Aggressive	No-suppressors	1.0	5.0	1.97	0.17	1.75	1.83	2.03	2.09	2.14	0.02156

Table 6: Detailed multi-target results.

Regression Size	Merge Mode	Strategy	$\frac{HF_s}{HF_n} > 1.0$	count	mean	std. dev	min	25%	50%	75%	max	p-val
SpiderMonkey												
2	Round-robin	Half-swarm	0.74	123.0	1.35	0.56	0.11	0.99	1.39	1.66	3.26	0.0
2	Round-robin	Triggers-only	0.77	13.0	2.22	1.5	0.01	1.55	2.49	2.74	5.23	0.00537
2	Round-robin	No-suppressors	0.6	123.0	1.1	0.44	0.07	0.75	1.15	1.4	2.08	0.01657
2	Subsumption	Half-swarm	0.72	155.0	1.36	0.57	0.01	0.93	1.37	1.74	2.88	0.0
2	Subsumption	Triggers-only	0.82	17.0	2.31	1.35	0.01	2.05	2.5	2.72	5.18	0.0007
2	Subsumption	No-suppressors	0.59	155.0	1.12	0.45	0.0	0.78	1.13	1.49	2.08	0.00154
2	Aggressive	Half-swarm	0.8	124.0	1.69	0.77	0.1	1.18	1.64	2.15	3.41	0.0
2	Aggressive	Triggers-only	0.92	12.0	2.8	1.5	0.1	2.21	2.5	2.93	6.02	0.00185
2	Aggressive	No-suppressors	0.61	122.0	1.18	0.49	0.09	0.84	1.19	1.55	2.36	0.00033
3	Round-robin	Half-swarm	0.8	142.0	1.29	0.39	0.12	1.06	1.33	1.57	2.22	0.0
3	Round-robin	Triggers-only	0.78	27.0	1.88	1.19	0.02	1.39	1.69	2.68	4.19	0.00114
3	Round-robin	No-suppressors	0.65	142.0	1.13	0.34	0.07	0.94	1.12	1.34	1.96	0.0
3	Subsumption	Half-swarm	0.79	143.0	1.29	0.4	0.17	1.04	1.28	1.55	2.34	0.0
3	Subsumption	Triggers-only	0.79	28.0	1.87	1.15	0.01	1.42	1.72	2.59	4.19	0.00072
3	Subsumption	No-suppressors	0.67	143.0	1.14	0.34	0.1	0.94	1.15	1.36	1.86	0.0
3	Aggressive	Half-swarm	0.87	142.0	1.63	0.58	0.09	1.25	1.6	2.07	3.23	0.0
3	Aggressive	Triggers-only	0.94	32.0	2.92	1.38	0.08	1.64	3.26	3.94	5.05	0.0
3	Aggressive	No-suppressors	0.73	142.0	1.23	0.4	0.13	0.98	1.22	1.53	2.21	0.0
4	Round-robin	Half-swarm	0.71	140.0	1.17	0.35	0.34	0.94	1.19	1.37	2.28	0.0
4	Round-robin	Triggers-only	0.94	33.0	1.86	0.8	0.09	1.51	1.74	2.12	3.92	0.0
4	Round-robin	No-suppressors	0.6	140.0	1.07	0.32	0.31	0.85	1.07	1.29	1.95	0.004
4	Subsumption	Half-swarm	0.69	141.0	1.18	0.37	0.36	0.93	1.17	1.4	2.5	0.0
4	Subsumption	Triggers-only	0.94	33.0	1.85	0.81	0.1	1.42	1.8	2.1	4.1	0.0
4	Subsumption	No-suppressors	0.6	140.0	1.06	0.32	0.32	0.83	1.07	1.27	1.77	0.00743
4	Aggressive	Half-swarm	0.82	140.0	1.5	0.55	0.36	1.15	1.46	1.78	3.23	0.0
4	Aggressive	Triggers-only	0.94	36.0	2.48	1.22	0.06	1.74	2.17	3.25	5.68	0.0
4	Aggressive	No-suppressors	0.67	140.0	1.17	0.38	0.29	0.9	1.16	1.41	2.2	0.0
5	Round-robin	Half-swarm	0.75	137.0	1.21	0.32	0.48	1.03	1.2	1.43	1.95	0.0
5	Round-robin	Triggers-only	0.71	38.0	1.84	1.13	0.01	0.97	1.79	2.5	4.38	7e-05
5	Round-robin	No-suppressors	0.66	136.0	1.12	0.29	0.45	0.92	1.13	1.35	1.79	2e-05
5	Subsumption	Half-swarm	0.76	140.0	1.22	0.32	0.5	1.01	1.22	1.44	1.96	0.0
5	Subsumption	Triggers-only	0.71	38.0	1.84	1.13	0.03	0.96	1.75	2.49	4.33	7e-05
5	Subsumption	No-suppressors	0.67	140.0	1.12	0.29	0.45	0.94	1.14	1.32	1.73	1e-05
5	Aggressive	Half-swarm	0.87	138.0	1.54	0.51	0.45	1.18	1.53	1.85	3.01	0.0
5	Aggressive	Triggers-only	0.85	34.0	2.53	1.22	0.19	1.46	2.74	3.42	4.86	0.0
5	Aggressive	No-suppressors	0.75	137.0	1.23	0.36	0.39	1.02	1.21	1.47	2.17	0.0
10	Round-robin	Half-swarm	0.73	145.0	1.14	0.21	0.58	0.99	1.14	1.27	1.78	0.0
10	Round-robin	Triggers-only	0.64	67.0	1.23	0.78	0.01	0.65	1.17	1.65	4.04	0.0309
10	Round-robin	No-suppressors	0.71	145.0	1.1	0.2	0.51	0.96	1.11	1.23	1.6	0.0
10	Subsumption	Half-swarm	0.7	145.0	1.13	0.21	0.54	0.99	1.12	1.26	1.92	0.0
10	Subsumption	Triggers-only	0.64	67.0	1.22	0.77	0.01	0.65	1.17	1.64	4.04	0.03178
10	Subsumption	No-suppressors	0.66	145.0	1.1	0.21	0.51	0.96	1.1	1.25	1.58	0.0
10	Aggressive	Half-swarm	0.92	145.0	1.44	0.32	0.63	1.25	1.45	1.65	2.4	0.0
10	Aggressive	Triggers-only	0.9	60.0	2.12	0.84	0.46	1.56	2.01	2.73	4.06	0.0
10	Aggressive	No-suppressors	0.81	145.0	1.24	0.26	0.5	1.07	1.25	1.39	1.99	0.0
20	Round-robin	Half-swarm	0.76	138.0	1.1	0.14	0.72	1.02	1.1	1.19	1.47	0.0
20	Round-robin	Triggers-only	0.47	91.0	1.04	0.62	0.01	0.54	0.93	1.45	2.45	0.62274
20	Round-robin	No-suppressors	0.75	138.0	1.09	0.15	0.68	1.0	1.1	1.2	1.44	0.0
20	Subsumption	Half-swarm	0.78	138.0	1.1	0.13	0.7	1.01	1.11	1.18	1.49	0.0
20	Subsumption	Triggers-only	0.49	90.0	1.04	0.62	0.0	0.59	0.97	1.46	2.44	0.57736
20	Subsumption	No-suppressors	0.76	138.0	1.1	0.14	0.74	1.01	1.08	1.2	1.45	0.0
20	Aggressive	Half-swarm	0.98	138.0	1.42	0.25	0.78	1.26	1.42	1.57	2.3	0.0
20	Aggressive	Triggers-only	0.93	96.0	1.92	0.62	0.28	1.48	1.94	2.29	3.37	0.0
20	Aggressive	No-suppressors	0.93	138.0	1.27	0.2	0.76	1.13	1.27	1.4	1.92	0.0

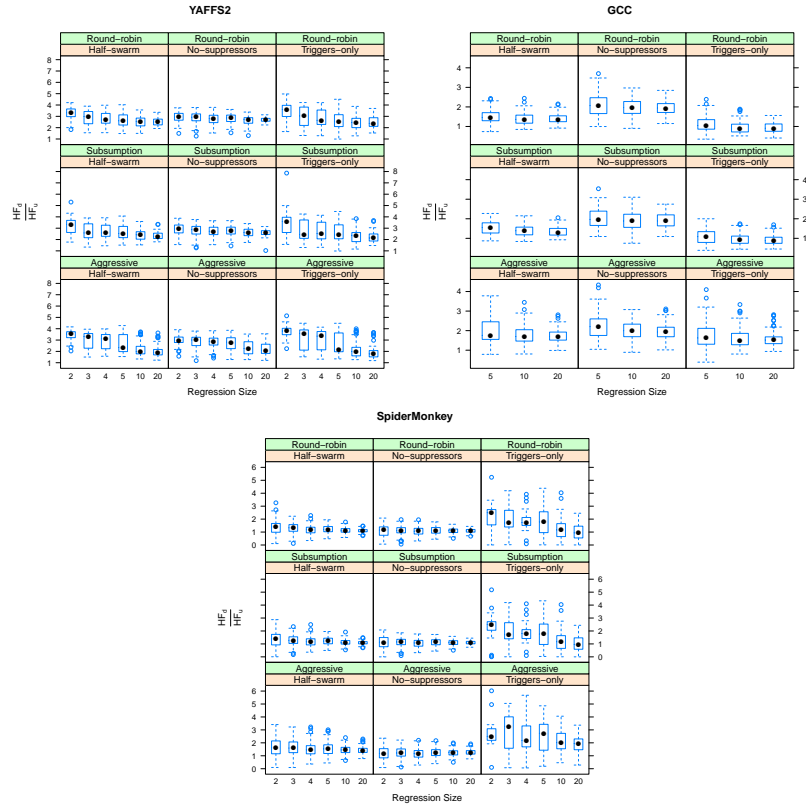


Figure 5: Multi-target strategies compared.

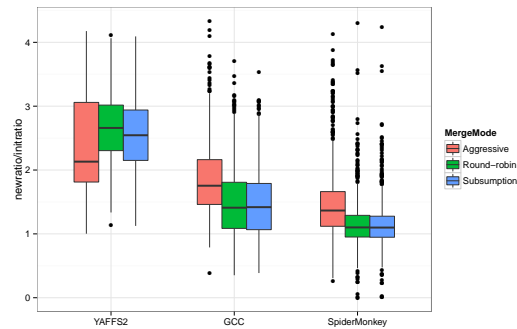


Figure 6: Merge strategies over all single-target strategies.

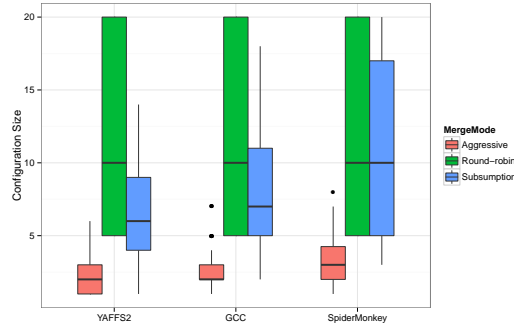


Figure 7: Number of targets after merging, by merge strategy.

Table 7: Detection rate of actual faults in the test suites generated by each technique in a 30-minute test suite generated by each test strategy.

Test Strategy	Fault #1	Fault #2	Fault #3	Fault #4	Fault #5
Undirected Swarm	13.6	0.07	0.24	0.26	0.07
Round-robin Half-swarm	31.9	0.19	0.35	0.56	0.29
Round-robin No-suppressors	34.2	0.26	0.17	0.46	0.69
Subsumption Half-swarm	33.0	0.24	0.12	0.10	0.29
Subsumption No-suppressors	33.1	0.31	0.29	0.31	0.46

experiments are based on 7 randomly chosen known (fixed) SpiderMonkey faults. For each of these faults, we targeted the statements in the code commit that introduced the fault. Evaluation was based on comparing 30 minute undirected and directed swarm suites, and counting how many times the fault was detected, on average, over a large (≥ 50) number of trials. For two of the faults, neither directed nor undirected testing ever detected the fault. The commit sizes for the remaining 5 faults were 40, 13, 5, 17, and 15 statements, respectively. Table 7 shows detection rates for undirected swarm testing and directed strategies, with the best detection rate for each fault in bold. Triggers-only is omitted from results, due to its difficulties producing valid SpiderMonkey tests, and Aggressive merging did not actually produce any additional merges over those provided by Subsumption. While no single strategy dominated all others, some basic points are clear: first, undirected swarm never had the best detection rate, and had the worst detection rate for 3 of the 5 faults. Second, Round-robin Half-swarm never had the worst detection rate, and had the best detection rate for 2 of the 5 faults, and improved the detection rate compared to undirected swarm testing by an average of 2.56x. Subsumption No-suppressors also always improved on undirected testing. Due to the large number of similar results (most runs did not detect a fault), however, these differences were only statistically significant for Fault #5.

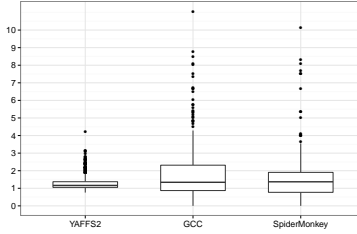


Figure 8: $\frac{HF_d}{HF_r}$, random vs. directed.

Observation 4: Directed random testing, for 5 real SpiderMonkey faults, could trigger real faults much more frequently than undirected testing. Round-robin Half-swarm was arguably the most effective approach.

6.4 RQ6: Comparison with Random Testing

We chose to use random testing without swarm configuration as an external evaluation. The reasons are twofold: (1) we could not find tools from other techniques (search based, concolic) that scale to our subject programs, and (2) if directed testing cannot improve the hitting fraction of targets over pure random testing, the applicability of our technique would be questionable. We compared the hitting fraction of targets in a single-targeted experiment (HF_d), with the hitting fraction in the test suites generated by traditional random testing (HF_r). Figure 8 illustrates the results. We used a paired t-test between HF_d and HF_r . Table 8 summarizes the results.

Table 8: The result of t-test comparing the hitting fraction of targets in directed random testing and random testing without swarm.

SUT	HF_r	HF_d	confidence-interval	p-value
YAFFS2	0.671	0.819	(0.126,0.170)	0.0000
GCC	0.342	0.425	(0.053,0.113)	0.0000
SpiderMonkey	0.198	0.276	(0.059,0.097)	0.0000

Observation 5: Directed swarm testing significantly ($p < 0.01$) increases hitting fractions over pure random testing.

7 Threats to Validity

Threats Due to Sampling Bias: Our results are based on results from three large open source software programs. While we believe that these programs are well tested examples of real-world programs, there is a possibility that they are not representative

of all software. Secondly, all our subjects are implemented in C, which may make the results not fully generalizable to other languages. Projects following other paradigms of development such as functional programming or object-oriented programming may also not be well represented.

Threats Due to Tools Used: Our results are dependent on the software we used to execute, collect and analyze our results. While we have tried to make sure that our results hold, using redundant checks over results, it is possible that we may have missed faults that may have an impact on our results.

Limited External Evaluation: We used pure random testing as our external evaluation. We are aware that there are other techniques that aim to cover particular targets code, most notably search-based techniques and symbolic execution. Unfortunately, the most versatile search based test generation tool we are aware of, EvoSuite [9]) only supports API-based test generation in Java. Only one of our test subjects (YAFFS2) is API-based, and it is written in C, using a complex feedback-based approach to method call choice that is based on file-system semantics [15]. For the other two subjects (GCC and SpiderMonkey) we were unable to find any applicable search-based or symbolic execution tools. At present, automated testing for production compilers is in practice largely performed using customized random testers. We further discuss our choice of subjects below.

8 Discussion

8.1 Subject Selection

Our subjects are chosen with two criteria in mind: first, they represent different kinds of features for swarm testing. YAFFS2 features are API calls, but (unlike the Java libraries more commonly used in the literature of API-call test generation), the calls modify a single, very complex program state (the file system itself) with complex dependencies. Features for SpiderMonkey testing using `jsfunfuzz` are actual production rules in a recursive generator, very difficult for a human engineer to understand (but easy to implement in a swarm tester). The complex recursive generation makes it an interesting subject to gauge the limits of our technique. Finally, test features in Csmith [30] are high-level semantic features of C programs, some of which do not correspond to simple grammar productions, and the features were devised to help compiler engineers deal with compilers with limited support for various C features, not for use in swarm testing. The second reason for selecting these subjects is that generating regression tests for subjects suitable for testing with current symbolic execution tools or search-based tools is not, we suspect, very difficult: existing techniques should suffice. However, many of the most critical software systems with complex input semantics are at this time only subjected to automated testing using custom-built random testers such as those we consider in this paper.

8.2 The Complexities of Configuration

While our results generally support the effectiveness of directed swarm testing, it is surprising how difficult it is to identify a single best strategy for directed swarm testing. Triggers-only is likely ineffective, but choosing between Half-Swarm, No-Suppressors, Round-Robin, Subsumption, and Aggressive strategies is not simple. In part we attribute this to the underlying complexity of what is happening in (directed) swarm testing: each configuration defines a (usually effectively infinite) set of tests. This is, of course, the point of random testing, that an unbounded number of diverse tests can be generated, using all available testing budget.

Swarm testing improves random testing in many cases, in the long run, by increasing the diversity of generated tests. This diversity can come with a price, however: for a fixed testing budget, because swarm testing improves diversity, the hitting fraction for many individual targets will be lower than for pure random testing (when swarm testing increases overall coverage, this is almost required — hitting more targets means hitting each target less often [17]). In fact, we noticed that comparing hitting fractions for undirected swarm testing and pure random testing, we often saw better hitting fractions for pure random testing, despite the fact that fault detection and overall coverage tend to show swarm testing performing much better for reasonably-long test runs [17]. Configuration strategy not only determines individual test behavior, but determines how quickly coverage saturates due to (lack of) diversity of tests created. Swarm testing produces very diverse tests; random testing without swarm configuration produces much less diverse tests. Our directed swarm testing strategies introduce a large number of choices in between these extremes, with a given focus [12]. Our experiments show that a variety of configuration methods can improve hitting fractions, but understanding how to best choose a strategy for, e.g., short vs. long budget directed testing is an open question we would like to address. However, the primary aim of directed regression tests is to catch faults quickly. One reasonable approach is to extend the diversity-centric ideas of swarm testing to strategy selection, and run in parallel directed tests for a change set using all of the viable strategies (e.g., all but Triggers-only).

9 Related Work

The most closely related work is our own previous work introducing swarm testing [17] and introducing the notions of triggers and suppressors [16].

There are several approaches for generating a test case that covers a chosen source code target. Of these, search-based testing [19, 23] and (dynamic) symbolic execution [11, 29] are the most notable ones. Symbolic execution [21] formulates an execution path in the program as a constraint formula problem and generates inputs that satisfy the path conditions and hence cover the target. Dynamic symbolic execution improves the scalability of pure symbolic execution by using information from concrete executions to replace over-complex constraints, simplifying problems of handling, e.g. system calls and pointers [11]. Search-based testing reduces the problem of covering a particular entity in the program into a search problem and uses search techniques, such as genetic algorithms and hill climbing, to solve this problem [19, 23].

There are many previous efforts to improve the test cases generated by random testing. Randoop [25] generates tests for object-oriented programs by calling random APIs, but uses feedback to guide test sequence creation. Nighthawk [2] uses genetic algorithms on top of a random tester to modify the configuration of the random tester to optimize it for a given goal (i.e., fitness function). Adaptive random testing [5,6] aims to improve random testing by using a distance measure to select more uniformly distributed tests, though its actual effectiveness in practice has been criticized [3]. ABP-based testing uses reinforcement learning to guide test generation [14].

To our knowledge, none of these approaches are applicable to the problem we address. While symbolic execution and search-based testing may be helpful for producing tests targeting a given element in source code, they are not always easy to apply to a target (such as a production quality compiler that takes as input full programs in a complex language), and symbolic execution in particular is often far less efficient than random testing [31]. Symbolic execution also often simply fails to scale to very large systems with complex input. The approach proposed in this paper is often trivial to apply to existing random test generators for complex software systems and, like pure random testing, has extremely low overhead (collecting coverage information on some random test runs is the only real cost, and this is only paid during data collection, not during new testing runs). While other methods may be suitable for generating targeted unit tests for concentrated code changes, we know of no other method that scales targeted automated testing to changes scattered throughout a large code base such as a compiler.

10 Conclusions and Future Work

In this paper we demonstrate that using collected statistics on code coverage and swarm testing, it is possible to produce *random regression tests* — truly random tests that nonetheless target specific source code targets. While results for the various strategies for *directed swarm testing* vary, in general the method is able to increase the frequency with which tests cover targeted code by a factor often more than 2x, and sometimes up to 8 or 9x. The most important aspect of this approach, unlike other methods of producing tests targeting certain coverage (symbolic or search-based approaches) is that it is readily applicable to existing, industrial-strength random testing tools for critical systems software, and therefore out-of-the-box scalable to applications such as testing production compilers and file systems. The changes in effectiveness of directed swarm testing, depending on the strategy chosen for balancing focusing testing and maintaining test case diversity, show the difficulty of understanding complex testing systems.

There are likely applications of directed swarm testing in addition to targeting just-changed code. For example, if static analysis indicates that a source code line may have a bug, but the analysis technique is subject to false positives, it may be useful to subject such lines to further scrutiny with targeted tests. Targeting source code that is very infrequently covered during extensive random testing, but covered enough to provide a basis for statistical estimation of triggers and suppressors may lead to covering code that the seldom-covered code dominates in the CFG, improving the overall effectiveness of large-scale random testing. Targeting faults, rather than source code

lines, can help improve suites for fault localization, by producing more failing tests to analyze. We believe there may be further practical applications of the combination of test suite statistics and variation in test case configurations.

References

- [1] ALIPOUR, M. A., AND GROCE, A. Bounded model checking and feature omission diversity. In *International Workshop on Constraints in Formal Verification* (2011).
- [2] ANDREWS, J. H., LI, F. C. H., AND MENZIES, T. Nighthawk: A two-level genetic-random unit test data generator. In *Proceedings of the Twenty-second IEEE/ACM International Conference on Automated Software Engineering* (New York, NY, USA, 2007), ASE '07, ACM, pp. 144–153.
- [3] ARCURI, A., AND BRIAND, L. Adaptive random testing: An illusion of effectiveness. In *International Symposium on Software Testing and Analysis* (2011), pp. 265–275.
- [4] ARCURI, A., IQBAL, M. Z. Z., AND BRIAND, L. C. Formal analysis of the effectiveness and predictability of random testing. In *International Symposium on Software Testing and Analysis* (2010), pp. 219–230.
- [5] CHEN, T. Y., KUO, F.-C., MERKEL, R. G., AND TSE, T. H. Adaptive random testing: The art of test case diversity. *J. Syst. Softw.* 83, 1 (Jan. 2010), 60–66.
- [6] CHEN, T. Y., LEUNG, H., AND MAK, I. Adaptive random testing. In *Advances in Computer Science-ASIAN 2004. Higher-Level Decision Making*. Springer, 2005, pp. 320–329.
- [7] CUOQ, P., MONATE, B., PACALET, A., PREVOSTO, V., REGEHR, J., YAKOBOWSKI, B., AND YANG, X. Testing static analyzers with randomly generated programs. In *NASA Formal Methods Symposium* (2012), pp. 120–125.
- [8] DE MOURA, L. M., AND BJØRNER, N. Z3: an efficient SMT solver. In *Tools and Algorithms for the Construction and Analysis of Systems* (2008), pp. 337–340.
- [9] FRASER, G., AND ARCURI, A. Evosuite: Automatic test suite generation for object-oriented software. In *Proceedings of the 19th ACM SIGSOFT Symposium and the 13th European Conference on Foundations of Software Engineering* (New York, NY, USA, 2011), ESEC/FSE '11, ACM, pp. 416–419.
- [10] GODEFROID, P., KIEZUN, A., AND LEVIN, M. Y. Grammar-based whitebox fuzzing. In *Proceedings of the 29th ACM SIGPLAN Conference on Programming Language Design and Implementation* (New York, NY, USA, 2008), PLDI '08, ACM, pp. 206–215.

- [11] GODEFROID, P., KLARLUND, N., AND SEN, K. Dart: Directed automated random testing. In *Proceedings of the 2005 ACM SIGPLAN Conference on Programming Language Design and Implementation* (New York, NY, USA, 2005), PLDI '05, ACM, pp. 213–223.
- [12] GROCE, A. (Quickly) testing the tester via path coverage. In *Workshop on Dynamic Analysis* (2009).
- [13] GROCE, A., ALIPOUR, M. A., ZHANG, C., CHEN, Y., AND REGEHR, J. Cause reduction for quick testing. In *Software Testing, Verification and Validation (ICST), 2014 IEEE Seventh International Conference on* (2014), IEEE, pp. 243–252.
- [14] GROCE, A., FERN, A., PINTO, J., BAUER, T., ALIPOUR, A., ERWIG, M., AND LOPEZ, C. Lightweight automated testing with adaptation-based programming. In *IEEE International Symposium on Software Reliability Engineering* (2012), pp. 161–170.
- [15] GROCE, A., HOLZMANN, G., AND JOSHI, R. Randomized differential testing as a prelude to formal verification. In *International Conference on Software Engineering* (2007), pp. 621–631.
- [16] GROCE, A., ZHANG, C., ALIPOUR, M., EIDE, E., CHEN, Y., AND REGEHR, J. Help, help, I'm being suppressed; The significance of suppressors in software testing. In *2013 IEEE 24th International Symposium on Software Reliability Engineering (ISSRE)* (Nov 2013), pp. 390–399.
- [17] GROCE, A., ZHANG, C., EIDE, E., CHEN, Y., AND REGEHR, J. Swarm testing. In *Proceedings of the 2012 International Symposium on Software Testing and Analysis* (New York, NY, USA, 2012), ISSTA 2012, ACM, pp. 78–88.
- [18] HAMLET, R. Random testing. In *Encyclopedia of Software Engineering*. Wiley, 1994, pp. 970–978.
- [19] HARMAN, M., MANSOURI, S. A., AND ZHANG, Y. Search-based software engineering: Trends, techniques and applications. *ACM Comput. Surv.* 45, 1 (Dec. 2012), 11:1–11:61.
- [20] KIM, S., WHITEHEAD, E., AND ZHANG, Y. Classifying software changes: Clean or buggy? *Software Engineering, IEEE Transactions on* 34, 2 (March 2008), 181–196.
- [21] KING, J. C. Symbolic execution and program testing. *Communications of the ACM* 19, 7 (1976), 385–394.
- [22] LE, V., AFSHARI, M., AND SU, Z. Compiler validation via equivalence modulo inputs. In *ACM SIGPLAN Conference on Programming Language Design and Implementation* (2014), pp. 216–226.

- [23] MCMINN, P. Search-based software testing: Past, present and future. In *Software Testing, Verification and Validation Workshops (ICSTW), 2011 IEEE Fourth International Conference on* (March 2011), pp. 153–163.
- [24] NAGAI, E., HASHIMOTO, A., AND ISHURA, N. Scaling up size and number of expressions in random testing of arithmetic optimization in c compilers. In *Workshop on Synthesis and System Integration of Mixed Information Technologies* (2013), pp. 88–93.
- [25] PACHECO, C., LAHIRI, S. K., ERNST, M. D., AND BALL, T. Feedback-directed random test generation. In *Proceedings of the 29th International Conference on Software Engineering* (Washington, DC, USA, 2007), ICSE '07, IEEE Computer Society, pp. 75–84.
- [26] RUDERMAN, J. Introducing jsfunfuzz. <https://www.squarefree.com/2007/08/02/introducing-jsfunfuzz/>.
- [27] VYUKOV, D. gosmith: Random Go program generator. <https://code.google.com/p/gosmith/>.
- [28] WILSON, E. B. Probable inference, the law of succession, and statistical inference. *J. of the American Statistical Assoc.* 22 (1927), 209–212.
- [29] XIE, T., TILLMANN, N., DE HALLEUX, J., AND SCHULTE, W. Fitness-guided path exploration in dynamic symbolic execution. In *Dependable Systems & Networks, 2009. DSN'09. IEEE/IFIP International Conference on* (2009), IEEE, pp. 359–368.
- [30] YANG, X., CHEN, Y., EIDE, E., AND REGEHR, J. Finding and understanding bugs in c compilers. In *Proceedings of the 32Nd ACM SIGPLAN Conference on Programming Language Design and Implementation* (New York, NY, USA, 2011), PLDI '11, ACM, pp. 283–294.
- [31] ZHANG, C., GROCE, A., AND ALIPOUR, M. A. Using test case reduction and prioritization to improve symbolic execution. In *International Symposium on Software Testing and Analysis* (2014), pp. 160–170.