

Using Symbolic Evaluation to Understand Behavior in Configurable Software Systems

Elnatan Reisner, Charles Song, Kin-Keung Ma, Jeffrey S. Foster, and Adam Porter
Computer Science Department, University of Maryland, College Park
{elnatan,csfalcon,kkma,jfoster,aporter}@cs.umd.edu

ABSTRACT

Many modern software systems are designed to be highly configurable, which increases flexibility but can make programs hard to test, analyze, and understand. We present an initial empirical study of how configuration options affect program behavior. We conjecture that, at certain levels of abstraction, configuration spaces are far smaller than the worst case, in which every configuration is distinct. We evaluated our conjecture by studying three configurable software systems: vsftpd, ngIRCd, and grep. We used symbolic evaluation to discover how the settings of run-time configuration options affect line, basic block, edge, and condition coverage for our subjects under a given test suite. Our results strongly suggest that for these subject programs, test suites, and configuration options, when abstracted in terms of the four coverage criteria above, configuration spaces are in fact much smaller than combinatorics would suggest and are effectively the composition of many small, self-contained groupings of options.

Categories and Subject Descriptors

D.2.5 [Software Engineering]: Testing and Debugging—*Symbolic execution, Testing tools*

General Terms

Measurement

Keywords

Empirical software engineering, software configurations, software testing and analysis

1. INTRODUCTION

Many modern software systems include numerous user-configurable options. For example, network servers typically let users configure the active port, the maximum number of connections, what commands are available, and so on.

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

ICSE '10, May 2-8 2010, Cape Town, South Africa
Copyright 2010 ACM 978-1-60558-719-6/10/05 ...\$10.00.

While this flexibility helps make software systems extensible, portable, and achieve good quality of service, it can also generate a huge number of configurations—in the worst case, every combination of option settings is a distinct configuration. This *software configuration space explosion* presents real challenges to software developers. It makes testing even more costly, as it significantly magnifies testing obligations; it makes static analysis much more difficult, as different configurations can be conflated together; and it generally complicates program understanding tasks.

In this paper, we present an initial empirical study exploring how various program behaviors change in relation to system configuration. We conjecture that at certain levels of abstraction, the software configuration space is much smaller than combinatorics might suggest. For example, consider a web server that can be configured with either sequential or concurrent connections, and to enable or disable PHP scripts. Since these two option settings likely do not affect each other, we probably only need two configurations (say, sequential with PHP and concurrent without PHP) to achieve full block coverage. If configuration options commonly follow this pattern, then in future work, new techniques and heuristics might be created to partition configuration spaces in ways that greatly simplify testing, analysis, and program understanding. As far as we are aware, our work is the first to formalize and precisely quantify the structure of configuration spaces. Previous understanding of configuration spaces is either anecdotal or based on indirect measurements (more discussion in Section 6).

To evaluate our conjecture, we studied three configurable subject systems: vsftpd, ngIRCd, and grep. For each system, we identified a sizable number of run-time configuration options to analyze, determined their possible settings, and created a test suite. We then ran the test suites using Otter, a symbolic evaluator [9, 7, 2] we developed.

In our study, we marked the selected configuration options as *symbolic*, meaning they represent unknowns that can take on any value. As Otter evaluates a program, if it encounters a branch that depends on a symbolic value, it conceptually forks execution and explores both possible branches. In this way, we used Otter to compute *all* possible program paths for *all* possible settings of the selected configuration options. This required tens or hundreds of thousands of runs, varying with the application, but that is a small fraction of the tens of millions or more runs that would have been needed had we naively enumerated all configurations.

We next projected the runs onto four types of structural coverage—line, basic block, edge, and condition coverage—

and used the results to find *interactions* among configuration options. We define an interaction to be a partial setting of configuration options such that specific coverage is *guaranteed* to occur under that setting, but is not guaranteed by any of its subsets. For example, if a and b are options, then a=0, b=1 is an interaction if it guarantees some coverage that setting a=0 or b=1 individually do not guarantee. Interactions are interesting because they define small subsets of the configuration space that provide meaningful additional coverage.

We computed interactions incrementally starting with combinations of zero options (i.e., coverage that occurs in all runs), then one option, then two, and so on. We continued until the accumulated guaranteed coverage equaled the maximum possible coverage across all runs. We found that for our subject systems and test suites, the largest interactions included between five and seven options. This is much smaller than the total number of options in any of the systems. Similarly, the total number of interactions is much smaller than what we would expect if we simply multiplied out the possible option combinations naively. These trends, and the others reported below, were essentially the same under all four coverage metrics.

Next, we looked at which specific interactions are needed to achieve high coverage. We found that most coverage is supplied by low-strength interactions, though all three programs had a few (one to three) *enabling options* that must be set a certain way to get maximum coverage. We also used a (non-optimal) greedy algorithm to pack together interactions into full configurations, in order to find the smallest set of configurations that would achieve full coverage. For example, interactions a=0, b=1 and c=0 can be joined into a single configuration, whereas a=0 and a=1 must go into different configurations. We found that we needed at most 9 configurations for vsftpd and for ngIRCd and 10 for grep to achieve full coverage. This suggests that for the programs and test suites used, the behavior of all configurations, when abstracted onto our coverage criteria, can be derived from the composition of a small number of interactions.

Finally, we created graphs showing all option interactions to better understand what allows us to achieve coverage with so few configurations. From this data and the previous analyses, we observe that for our programs, coverage metrics, test suites, and configuration spaces, many options do not interact with each other; that when they do interact, they often do so at low-strength; and that the interactions that exist often cluster into distinct groupings that can be combined into larger configurations.

In summary, our results strongly support our main conjecture: that in practical systems, when abstracting to specific program execution behaviors, the configuration space behaves less like a monolithic cross product of all option settings, and more like the union of smaller configuration spaces. We believe our work provides a basic but important starting point for understanding software configurability and for creating techniques and heuristics for scaling many software development tasks across large configuration spaces.

2. CONFIGURABLE SOFTWARE SYSTEMS

For this work, a configurable system is a generic code base and a set of mechanisms for implementing pre-planned variations in the code base’s structure and behavior. In practice, these variations are wide-ranging, covering hardware

```

1 ... else if (tunable_pasv_enable &&
2     str_equal_text(&p_sess->ftp_cmd_str, "EPSV"))
3 {
4     handle_pasv(p_sess, 1);
5 }
6 ... else if (tunable_write_enable &&
7     (tunable_anon_mkdir_write_enable ||
8     !p_sess->is_anonymous) &&
9     (str_equal_text(&p_sess->ftp_cmd_str, "MKD") ||
10    str_equal_text(&p_sess->ftp_cmd_str, "XMKD")))
11 handle_mkd(p_sess);
12 }

```

(a) Boolean configuration options (vsftpd)

```

13 if ((Conf_MaxJoins > 0) &&
14     (Channel_CountForUser(Client) >= Conf_MaxJoins))
15     return IRC_WriteStrClient(Client,
16                             ERR_TOOMANYCHANNELS_MSG,
17                             Client_ID(Client), channame);

```

(b) Integer-valued configuration options (ngIRCd)

```

18 else if(Conf_OperCanMode) {
19     /* IRC—Operators can use MODE as well */
20     if (Client_OperByMe(Origin)) {
21         modeok = true;
22         if (Conf_OperServerMode)
23             use_servermode = true; /* Change Origin to Server */
24     }
25 }
26 ...
27 if (use_servermode)
28     Origin = Client_ThisServer();

```

(c) Nested conditionals (ngIRCd)

```

29 not_text =
30 ((binary_files == BINARY_BINARY_FILES && !out_quiet)
31  || binary_files == WITHOUT_MATCH_BINARY_FILES)
32  && memchr (bufbeg, eol ? '\0' : '\200', buflim - bufbeg));
33 if (not_text &&
34     binary_files == WITHOUT_MATCH_BINARY_FILES)
35     return 0;
36 done_on_match += not_text;
37 out_quiet += not_text;

```

(d) Options being passed through the program (grep)

Figure 1: Uses of configuration options (bolded).

and operating system platforms (e.g., Windows vs. Linux), software versions (e.g., MySQL 5.0 vs. MySQL 5.1), runtime features (e.g., enable/disable debugging output), and others. In this paper, we focus on run-time configuration options, which are usually given values via configuration files or command-line parameters. A *configuration* is a mapping of configuration options to their settings.

Figure 1 illustrates several ways that run-time configuration options can be used, and explains why understanding their usage requires fairly sophisticated technology. All of these examples come from our experimental subject programs, which are written in C. In this figure, configuration options are shown in boldface.

The example in Figure 1(a) shows a section of vsftpd’s command loop, which receives a command and then uses a

long sequence of conditionals to interpret the command and carry out the appropriate action. The example shows two such conditionals that also depend on configuration options (all of which begin with `tunable_` in `vsftpd`). In this case, the configuration options enable certain commands, and the enabling condition can either be simply the current setting of the option (as on line 1) or may involve an interaction between multiple options (as on lines 6–7).

Not all options need to be booleans, of course. Figure 1(b) shows code from `ngIRCd` in which `Conf_MaxJoins` is an integer option that, if positive (line 13), gives the maximum number of channels a user can join (line 14). In this example, error processing occurs if the user tries to join too many channels.

Figure 1(c) shows a different example in which two configuration options are tested in nested conditionals. This illustrates that it is insufficient to look at tests of configuration options in isolation; we also need to understand how they may interact based on the program’s structure. Moreover, in this example, if both options are enabled then `use_servermode` is set on line 23, and its value is then tested on line 27. This shows that the values of configuration options can be indirectly carried through the state of the program.

Figure 1(d) shows another example of using configuration options indirectly. Here `not_text` is assigned the result of a complex test involving configuration options, and is then used in the conditional (lines 33–34) to change the current setting of two other configuration options (lines 36–37).

3. SYMBOLIC EVALUATION

To understand how configurations resemble and differ from each other, we have to capture their effect on a system’s runtime behavior. As we saw above, configuration options can be used in complex ways, and thus simple approaches such as searching through code for option names will be insufficient. Instead, we use symbolic evaluation [9] to capture all executions a program can take under any configuration.

Our symbolic evaluator, Otter,¹ is essentially a C source code interpreter, with one key difference. We allow the programmer to designate some values as *symbolic*, meaning they represent unknowns that may take on any value. Otter tracks these values as they flow through the program, and conceptually forks execution if a conditional depends on a symbolic value. Thus, if it runs to completion, Otter will simulate all paths through the program that are reachable for any values that the symbolic data can take.

To illustrate how Otter works, consider the example C source code in Figure 2(a). This program includes input variables `a`, `b`, `c`, `d`, and `input`. The first four are intended to represent run-time configuration options, and so we initialize them on lines 1–2 with *symbolic values* α , β , γ , and δ , respectively. (In the implementation, the content of a variable `v` is made symbolic with a special call `__SYMBOLIC(&v)`.) The last variable, `input`, represents program inputs other than configuration options. Thus we leave it concrete, and it must be supplied by the user (e.g., as part of `argv` (not shown)).

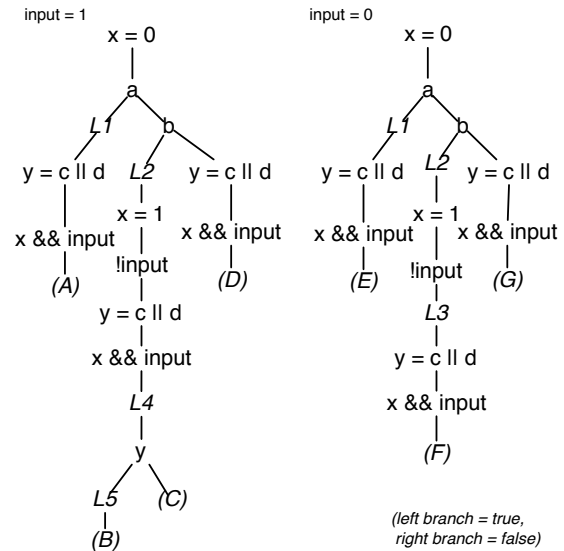
We have indicated five lines, represented by comments L1–L5, whose coverage we are interested in. Figure 2(b) shows the sets of paths explored by Otter as *execution trees*

```

1 int a= $\alpha$ , b= $\beta$ ,
2   c= $\gamma$ , d= $\delta$ ; // symbolic
3 int input=...; // concrete
4 int x = 0;
5 if (a) {
6   /* L1 */
7 } else if (b) {
8   /* L2 */
9   x = 1;
10  if (!input) {
11
12  }
13 }
14 int y = c || d;
15 if (x && input) {
16   /* L4 */
17   if (y) {
18     /* L5 */
19   }
20 }

```

(a) Example program



(b) Full execution trees

Figure 2: Example symbolic evaluation.

for two concrete “test cases” for this program: the tree for `input=1` is on the left, and the tree for `input=0` is on the right. Here nodes correspond to program statements, and branches represent places where Otter has a choice and hence “forks,” exploring both possible paths.

For example, consider the tree with `input=1`. All executions begin by setting `x` to 0 and then testing the value of `a`, which at this point contains α . Since there are no constraints on α , both branches are possible. For the sake of simplicity, we will assume that α and the other symbolic values may only represent 0 and 1, but Otter fully models symbolic integers as arbitrary 32-bit quantities.

Otter forks its execution at the test of `a`. First, it assumes $\alpha = 1$ and reaches L1 (left branch). It then falls through to line 14 (the assignment to `y`) and performs the test on line 15 (`x && input`). This test is false, since `x` was set to 0 earlier, hence Otter does not fork. We label this path through the execution tree as (A). Note that here we model `x && input` as a single test, rather than treating `&&` as short-circuiting. This is sound when the right-hand side of `&&` has no side effects, and provides more sensible coverage metrics; we discuss this more in Section 3.2.

Notice that as we explored path (A), we made some decisions about the settings of symbolic values, specifically that $\alpha = 1$. We call this and any other constraints placed on the

¹DART [7] and EXE [3] are two well known symbolic evaluators. By coincidence, Dart and Exe are the names of two rivers in Devon, England. The others are the Otter, the Tamar, the Taw, the Teign, and the Torridge.

symbolic values a *path condition*. Here, path (A) covers L1, and so any configuration that sets $a=1$ (corresponding to $\alpha = 1$), with arbitrary choices for β , γ , and δ , will cover L1. This is what makes symbolic evaluation so powerful: With a single predicate we characterized the behavior of many possible concrete choices of symbolic inputs.

Otter continues by returning to the last place it forked and trying to explore the other path. In this case, it returns to the conditional on line 5, adds $\alpha = 0$ to the path condition, and continues exploring the execution tree. Each time Otter encounters a conditional, it calls a Satisfiability Modulo Theory (SMT) solver to determine which branches (possibly both) of the conditional are possible based on the current path condition.

There are a few other interesting things to notice about these execution trees. First, consider the execution path labeled (B). Because we have assumed $\beta = 1$ on this path, we set $x=1$, and hence $x \ \&\& \ \text{input}$ is true, allowing us to reach L4 and L5. This is analogous to the example in Figure 1(c), in which a configuration option choice resulted in a change to the program state (setting $x=1$) that allowed us to cover some additional code. Also, notice that if $\text{input}=1$, there is no way to reach L3, no matter how we set the symbolic values. Hence coverage depends on choices of both symbolic values and concrete inputs.

In total, there are four paths that can be explored when $\text{input}=1$, and three paths when $\text{input}=0$. However, there are 2^5 possible assignments to the 5 input variables. Hence using symbolic evaluation for these test cases enables us to gather full coverage information with only 7 paths, rather than the 32 runs required if we had tried all possible combinations of symbolic and concrete inputs.

3.1 Guaranteed Coverage

Otter forms the basis for our empirical study: For each subject program, we selected a number of configuration options, marked them as symbolic, and then used Otter to execute a set of test cases. The resulting execution trees contain all possible paths executed under all configuration option settings for those test cases.

Without further analysis, these paths tell us only a little about our subject programs. By definition, each path explored for a particular test case is distinct from all the other paths for the same test case. Thus with no abstraction, every configuration option combination given by a path is unique. For example, in Figure 2(b), there are four distinct paths if $\text{input}=1$, representing four distinct settings of configuration options. Thus far, then, we only know that that is fewer than the 16 paths we might naively expect.

However, if we are interested in more abstract properties of the program, then paths are no longer unique, and the configuration space collapses further. For example, suppose we are only interested in covering L2 in Figure 2. Then we can see that paths (A) and (D) are irrelevant, and either path (B) or (C) is sufficient.

For our study, we project the symbolic evaluation results onto line, block, edge, and condition coverage. The principal tool we use to relate configuration options to coverage is *guaranteed coverage*.

DEFINITION 1. *Given a test suite and a coverage criterion, we say that a predicate p over the (initial settings of the) configuration options guarantees coverage of program entity X if there exists some test case in the test suite such*

that any configuration satisfying p is guaranteed to cover X .

For example, from Figure 2(b) we can see that any configuration satisfying $\alpha = 0 \wedge \beta = 1$ (i.e., $a=0, b=1$) is guaranteed to cover L2, no matter the choice of γ and δ .

We can use Otter’s output to compute the guaranteed coverage for a predicate p , which we will write as $Cov(p)$. We first find $Cov^T(p)$, the coverage guaranteed under p by test case T , for each test case; then, $Cov(p) = \bigcup_T Cov^T(p)$. To compute $Cov^T(p)$, let p_i^T be the path conditions from T ’s symbolic evaluation, and let $C^T(p_i^T)$ be the covered lines (or blocks, edges, conditions, etc.) that occur in that path. Then $Cov^T(p)$ is

$$\begin{aligned} Consistent^T(p) &= \{p_i^T \mid SAT(p_i^T \wedge p)\} \\ Cov^T(p) &= \bigcap_{q \in Consistent^T(p)} C^T(q) \end{aligned}$$

In words, first we compute the set of path conditions p_i^T such that p and p_i^T are consistent. If this holds for p_i^T , the items in $C^T(p_i^T)$ may be covered if p is true. Since our symbolic evaluator explores all possible program paths, the intersection of these sets for all such p_i^T is the set guaranteed to be covered if p is true.

Going back to Figure 2, here are some predicates and the coverage they guarantee given the test cases $\text{input}=1$ and $\text{input}=0$. We abbreviate $\alpha = 1$ as α and $\alpha = 0$ as $\neg\alpha$.

p	$Consistent(p)$ ($\text{input} = 1$)	$Consistent(p)$ ($\text{input} = 0$)	$Cov(p)$
α	(A)	(E)	{L1}
β	(A), (B), (C)	(E), (F)	\emptyset
$\neg\alpha$	(B), (C), (D)	(F), (G)	\emptyset
$\neg\alpha \wedge \beta$	(B), (C)	(F)	{L2, L3, L4}
$\neg\alpha \wedge \beta \wedge \gamma$	(B)	(F)	{L2, L3, L4, L5}

As we show in Section 5, we can use guaranteed coverage to discover *interactions* among options.

DEFINITION 2. *An interaction is a conjunction of option settings $S = \bigwedge_i (x_i = v_i)$ that guarantees coverage that is not guaranteed by any subset of (the conjuncts of) S .*

For example, $Cov(\alpha = 0 \wedge \beta = 1)$ is a strict superset of $Cov(\alpha = 0) \cup Cov(\beta = 1)$, so $\alpha = 0 \wedge \beta = 1$ is an interaction. Informally, interactions indicate option combinations that are “interesting” because they guarantee some new coverage.

DEFINITION 3. *The strength of an interaction is the number of option settings it contains.*

For example, $\alpha = 0 \wedge \beta = 1$ has strength 2. Lower-strength interactions place fewer requirements on configurations, whereas higher-strength interactions require more options to be set in particular ways to achieve their coverage.

3.2 Implementation

Otter is written in OCaml, and it uses CIL [11] as a front end to parse C programs and transform them into an easier-to-use intermediate representation.

The general approach used by Otter mimics KLEE [2]. A symbolic value in Otter represents a sequence of untyped bits, e.g., a 32-bit symbolic integer is treated as a vector with 32 symbolic bits in Otter. This low-level representation is important because many C programs perform bit manipulations that must be modeled accurately. When a symbolic expression has to be evaluated, Otter invokes STP [6], an SMT solver optimized for bit vectors and arrays.

Otter supports all the features of C we found necessary for our subject programs, including pointer arithmetic, arrays, function pointers, variadic functions, and type casts. Loops, which can cause symbolic evaluation to “get stuck” as it tries to unroll the loop an unbounded (or extremely large) number of times, were not a major obstacle in our study: configuration options almost never influenced loop bounds (see Section 4), so almost all loops were executed in the usual way, with the concrete test cases determining the number of loop iterations.

Otter currently does not handle dereferencing symbolic pointer values, floating-point arithmetic, or inline assembly. In all cases, these features either do not appear in our subject programs or do not affect the results of our study. Otter also does not support multiple processes, which do occur in vsftpd and ngIRCd. For vsftpd, in which `fork()` spawns a subprocess that handles client commands, we interpret `fork()` as the starting point of that subprocess and ignore the parent process, since it would simply cycle around a loop. For ngIRCd, where the child process parses an IP address and passes the result to the parent, we treat `fork()` as a branching point—we run both subprocesses, but we ignore the child process’s output, instead supplying the input expected by the parent process as part of the test case.

All of our subject programs interact with the operating system in some way. Thus, we developed “mock” libraries that simulate a file system, network, and other needed OS components. Our libraries also allow test cases to control the contents of files, data sent over the network, and so on. Our mock library functions are written in C and are executed by Otter just like any other program. For example, we simulate a file with a character array that also tracks the current position at which the file is to be read or written.

As Otter executes, it records the program paths explored so that we can later compute line, block, edge, and condition coverage. The precise definitions of these metrics demand some elaboration, because Otter runs on CIL’s representation of the input program, which is simplified to use only a subset of the full C language.

To compute line coverage, we record which CIL statements Otter executes and project that back to the original source lines using a mapping maintained by CIL.

For block and edge coverage, we group CIL statements into basic blocks, which are sequences of statements that start at a function entry or a join point; do not contain any join point after the first statement; end in a function call, `return`, `goto`, or conditional; or fall through to a join point. Normally, CIL expands short-circuiting logical operators `&&` and `||` into sequences of branches. However, for line, block, and edge coverage, we disable that expansion as long as the right operand has no side effect, so that both operands are computed in the same basic block. Then to compute block coverage, we record which basic blocks are executed, and to compute edge coverage, we compute which control-flow edges between basic blocks are traversed.

Lastly, for condition coverage, we enable expansion of `&&` and `||`, so that each part of a compound condition is always in its own basic block. We then compute how many conditions—that is, how many branches—are taken in the expanded program.

4. SUBJECT PROGRAMS

The subject programs for our study are vsftpd, a widely

	vsftpd	ngIRCd	grep
Version	2.0.7	0.12.0	2.4.2
# Lines (sloccount)	10,482	13,601	9,124
# Lines (executable)	4,112	4,387	3,302
# Basic blocks	4,584	6,742	5,033
# Edges	5,033	7,374	6,332
# Conditions	2,528	3,432	4,094
# Test cases	64	142	113
# Symbolic conf. opts.	30	13	18
Boolean	20	5	14
Integer	10	8	4
# Concrete conf. opts.	65	16	4
Full config, test space	1.4×10^{11}	4.2×10^7	6.7×10^7

Figure 3: Subject program statistics.

used secure FTP daemon; ngIRCd, the “next generation IRC daemon”; and GNU grep, a popular text search utility. All of our subject programs are written in C. Each has multiple configuration options that can be set either in system configuration files or through command-line parameters.

Figure 3 gives descriptive statistics for each subject program. The top two rows list the program version numbers and lines of code as computed by `sloccount`. The next group of rows lists the number of executable lines, basic blocks, edges, and conditions; these four metrics are what we measure code coverage against, and they are based on the CIL representation of the program, as discussed in Section 3.2. To get more accurate measurements, we removed some unreachable code before passing the sources to CIL. We eliminated out 4 unreachable functions in grep, and we eliminated 3 files in vsftpd that are reachable only in two-process mode, which we disabled because Otter does not support multiprocess symbolic evaluation.

One thing to note is that there are more basic blocks than executable lines of code in all 3 programs. This occurs because, in many cases, single lines form multiple blocks. For example, a line that contains a `for` loop will have at least two blocks (for the initializer and the guard), and lines with multiple function calls are broken into separate blocks according to our definition.

The next row in Figure 3 lists the number of test cases. In creating these test cases, we attempted both to cover the major functionality of the system and to maximize overall line coverage. We stopped creating new tests when we reached a point of diminishing returns for our efforts. For example, much of the code left uncovered by our test suites handles system call failures, such as `malloc()` returning `NULL`; modeling these failures would have greatly increased the number of execution paths (and hence analysis time) without shedding extra light on the configuration spaces of these programs. Other uncovered code would have required significantly extending Otter—e.g., to handle asynchronous signals—which was beyond of the scope of this initial study.

Vsftpd does not come with its own test suite, so we developed tests to exercise its major functionality such as logging in; listing, downloading, and renaming files; asking for system information; and handling invalid commands.

NgIRCd also does not come with its own test suite, so we created tests based on the IRC functionality defined in RFCs 1459, 2812 and 2813. Our tests cover most of the client-server commands (e.g., client registration, channel join/part, messaging and queries) and a few of the server-server commands (e.g., connection establishment, state exchange), with both valid and invalid inputs.

	vsftpd	ngIRCd	grep
# Paths			
Line, Block, Edge	30,304	53,205	625,181
Condition	136,320	95,637	764,201
Average # Paths			
Line, Block, Edge	474	375	5,533
Condition	2,130	674	6,763
Coverage			
Line	62%	73%	75%
Block	63%	66%	63%
Edge	56%	61%	58%
Condition	49%	57%	52%
# Examined opts/tot			
Line, Block, Edge	22/30	13/13	17/18
Condition	24/30	13/13	17/18

Figure 4: Summary of symbolic evaluation.

Grep comes with a test suite consisting of hundreds of tests. To build our test suite for this study, we ran all the test cases in Otter to determine their line coverage. Then, without sacrificing total line coverage, we selected 70 test cases from the original suite for our study. Next, we created 43 new test cases to improve overall line coverage. The final analysis was done using these 113 test cases.

Figure 3 next shows the counts of the configuration options. We give the total number of analyzed configuration options, i.e., those that we treated as symbolic, and also break them down by type (boolean or integer). We also list the number of configuration options we left as concrete. (For vsftpd, this count excludes 27 configuration options that were not present in the code after preprocessing.)

We decided to leave some options concrete based on two criteria: whether the option was likely to expose meaningful behaviors, and our desire to limit total analysis effort. We focused on boolean and integer options, and so we left more complex options (e.g., strings) concrete, with one exception: one of grep’s string options is only three-valued, so we considered it an integer option and set it symbolic.

For many integer options, only a few values (or classes of values, such as $x > 0$) are important in determining program execution, so Otter only had to consider these values. A few of the integer options, however, led to an unmanageable number of paths when made symbolic, e.g., because they were passed to `printf` or used as a loop bound. For these options, we either constrained them to range over a small number of values (chosen to maximize coverage), or left them concrete. Even then, the sheer number of symbolic options led to lengthy executions, so we continued reverting additional options to concrete values until the executions and subsequent analysis were manageable. This approach allowed us to run Otter numerous times on each program, to explore different scenarios, and to experiment with different kinds of analysis techniques. We used default values for the concrete configuration options, except the one used to force single-process mode in vsftpd.

Finally, we show how many executions would be required if we ran every test case under every possible configuration, given the number of distinct values each symbolic option could take. We will contrast these extremely large numbers with the results of symbolic evaluation in the next section.

5. DATA AND ANALYSIS

We ran our test suites in Otter, with symbolic configura-

tion options as discussed above. We then performed substantial analysis on the results to explore the configuration space of each subject program. To do this we used the Skoll system, developed and housed at UMD [12]. Skoll allows users to define configurable QA tasks and run them across large virtual computing grids. For this work we used around 40 client machines. The final results reported in this section required about two weeks of elapsed time.

Figure 4 summarizes Otter’s runs. The first two rows show the total number of paths executed by Otter, summed across all tests. This is dramatically smaller than the number of executions that would have been necessary had we instead naively run each test case under all possible configuration option combinations: 0.0001% of the naive count for vsftpd, 0.2% for ngIRCd, and 1% for grep. Also, recall that these are actually *all* possible paths for these test suites under any settings of the symbolic configuration options, given Otter’s simulated environment.

Notice that condition coverage, which has logical operators expanded into sequences of conditionals as discussed in Section 3.2, has many more paths. This effect is most pronounced for vsftpd, which more than quadruples the number of paths because it contains many logical expressions that test multiple configuration options at once. For example, `if (x||y||z)` would yield at most two paths before expansion, but four paths after.

To aid comparison across our subjects, we next show the total number of paths averaged over all test cases.

Figure 4 then lists the coverage achieved by these paths, i.e., the maximum coverage achievable for these test suites considering all possible configurations, except options left concrete. We manually inspected the uncovered lines and found that approximately another 10% of vsftpd and ngIRCd and 2% of grep comprise code for handling low-level errors, and another 11% of vsftpd (in addition to the three files we removed) is unreachable in single-process mode. If we adjust for the error handling and unreachable code, our test suites’ line coverage is near or exceeds 80% for all subject programs. Covering the remaining code would in many cases have required adding new mocked libraries, adding further symbolic configuration options, etc. Overall, however, based on our analysis of these systems, we believe that the test cases are reasonably comprehensive and are sufficient to expose much of the configurable behavior of the subject programs.

The next group of rows shows the number of configuration options that appear in at least one path condition (i.e., were constrained in at least one path and thus distinguish different execution paths) versus the total number of options set symbolic. In grep, the one unused option was only “used” when being printed, which did not affect any execution path. In vsftpd, there were six unused options. One case was similar to grep—a configuration option specified a port number, which is ignored by our mock system. Three other options could have been covered with additional tests; the remaining two options cannot be touched without changing the settings of some of the configurations options we left concrete.

Notice that Otter constrains two more options with condition coverage than under the other metrics. This occurs because of the expansion of logical operators into sequences of conditionals. For example, under line, block, and edge coverage, the condition `if (x||1)` would be treated as a single branch that Otter would treat as always true. But under condition coverage, the conditional would be expanded, and

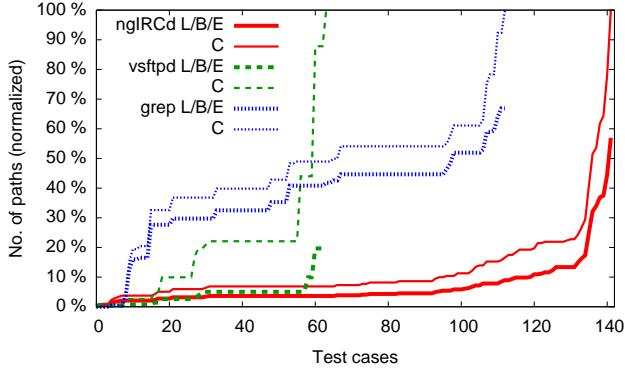


Figure 5: Number of paths per test case (L/B/E=line/block/edge, C=condition).

Otter would see if (x) first, causing it to branch on x .

Figure 5 plots the number of paths executed by each test case for each program, both with unexpanded logical operators (L/B/E) and expanded (C). The x -axis is sorted from the fewest to the most paths, and the y -axis is the percentage of paths relative to the highest number of paths seen in any test case for the expanded (C) version of the program.

One interesting feature of Figure 5 is that, for vsftpd and grep, the numbers of paths of different test cases cluster into a handful of groups (indicated by the plateaus in the graph). This suggests that within a group, the test cases branch on the configuration options in essentially the same manner (most likely because the programs employ common segments of code to test the configuration options). In ngIRCd, this clustering also appears but is less pronounced.

Finally, recall from Figure 3 that grep, despite still having many fewer paths than configurations, stands out as having a much larger number of paths than the other programs. We believe this is due to grep’s design. In runs of grep with valid inputs, most of grep’s code is executed. Therefore many of its configuration options will typically be used, resulting in significant branching in Otter. In contrast, many of vsftpd and ngIRCd’s options are not necessarily used in every run. This can be seen clearly in Figure 5: only a handful of vsftpd and ngIRCd’s tests exercise more than 25% of the paths, while only a handful of grep’s tests exercise fewer than 25%.

5.1 Interactions by Strength

Next, we used our guaranteed coverage analysis to explore which configuration option interactions (Section 3.1) are actually required to achieve the line, block, edge, and condition coverage reported in Figure 4. First, we computed $Cov(true)$, which we call *guaranteed 0-way coverage*. These are coverage elements that are guaranteed to be covered for any choice of options. Here when we refer to t -way coverage, t is the interaction strength. Then for every possible option setting $x = v$, we computed $Cov(x = v)$. The union of these sets is the *guaranteed 1-way coverage*, and it captures what coverage elements will definitely be covered by 1-way interactions. Next, we computed $Cov(x_1 = v_1 \wedge x_2 = v_2)$ for all possible pairs of option settings, which is *guaranteed 2-way coverage*. Similarly, we continue to increase the number of options in the interactions until $Cov(x_1 = v_1 \wedge x_2 = v_2 \wedge \dots)$ reaches the maximum possible coverage.

For boolean options, the possible settings are clearly 0

	$t=1$	$t=2$	$t=3$	$t=4$	$t=5$	$t=6$	$t=7$
vsftpd							
Line	7	4	3	16	5	6	2
Block	7	4	3	16	6	6	2
Edge	9	4	4	27	7	7	2
Condition	9	4	4	32	14	9	2
ngIRCd							
Line	11	19	33	117	144	111	-
Block	15	25	33	118	147	111	-
Edge	17	29	37	125	159	111	-
Condition	17	33	37	131	174	111	-
grep							
Line	13	27	36	7	5	-	-
Block	14	34	37	7	5	-	-
Edge	23	37	45	11	7	-	-
Condition	23	45	49	16	9	2	-

Figure 6: Number of interactions at each strength.

and 1. For integer-valued options that we constrained (as described in Section 4), we used those chosen values; for the remaining integer options, we solved the path conditions discovered by Otter and manually inspected the code to find appropriate concrete settings.

Figure 6 shows the number of interactions at each interaction strength. The first thing to notice is that the maximum interaction strength is always seven or less. This is significantly lower than the number of options in each program. We also see that the number of interactions is quite small relative to total number of interactions that are theoretically possible. For example, grep has 14 boolean options, which by themselves lead to $(14 \text{ choose } 2) \times 4 = 364$ possible 2-way interactions just with those options alone, yet we see at most 45 2-way interactions for grep.

Also notice that there is little variation across different coverage criteria—they have remarkably similar numbers of interactions. We investigated further and found the majority of interactions are actually identical across all four criteria. This is encouraging, because it indicates that many interactions are insensitive to the particular coverage criterion.

For ngIRCd, there are significantly more interactions at higher strength than for the other subject programs. This is because almost all of ngIRCd’s integer options can take on many different values across our test suite, magnifying the number of interactions.

Finally, we can see that the number of interactions peak around $t = 4$ for vsftpd, $t = 4$ or 5 for ngIRCd, and $t = 2$ or 3 for grep. We believe this corresponds to the number of *enabling options* in these programs, discussed more below.

5.2 Guaranteed Cov. by Interaction Strength

Figure 7 presents the interaction data in terms of coverage. The x -axis is the t -way interaction strength, and the y -axis is the percentage of the maximum possible coverage. Note that higher-level guaranteed coverage always includes the lower level, e.g., if a line is covered no matter what the settings are (0-way), then it is certainly covered under particular settings (1-way or higher). As it turns out, the trend lines for all four coverage criteria are essentially the same for a given program, and so the plot shows a region enclosing each data set. In ngIRCd, the only program with some slight variation, line coverage corresponds to the upper boundary of the region, and edge, block, and condition coverage to the lower boundary. This commonality across coverage criteria echoes the same trend we saw in Figure 6.

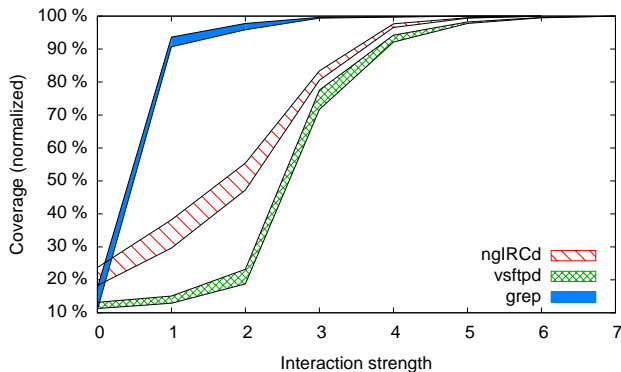


Figure 7: Guaranteed coverage versus interaction strength.

We also notice in this figure that the right-most portion of each region adds little to the overall coverage. Thus, for these programs and test suites, high-strength interactions are not needed to cover most of the code. We can also see that vsftpd gains coverage slowly but then spikes with 3-way interactions, and grep has a similar spike with 1-way interactions. This suggests the presence of *enabling options*, which must be set a certain way for the program to exhibit large parts of its behavior. For example, for vsftpd (in single-process mode), the enabling options must ensure local logins and SSL support are turned off, and anonymous logins are turned on. For grep, either grep or egrep mode must be enabled to reach most of grep’s code; fgrep mode touches little code. ngIRCd also has enabling options that account for the increasing coverage up to interaction strength three, but the effect of these options are less pronounced.

These enabling options also show up in Figure 6. For example, in that figure we can see that most of vsftpd’s interactions are strength $t = 4$ or greater, i.e., they generally involve the three enabling options plus additional options.

5.3 Minimal Covering Configuration Sets

Our results so far show that low-strength interactions can cover most of the code. Next, we investigated how interactions can be packed together to form complete configurations, which assign values to all configuration options. For example, 1-way interactions $a=0$ and $b=0$ are consistent and can be packed into the same configuration, but $a=0$ and $a=1$ are contradictory and must go in different configurations.

We developed a greedy algorithm that packs options together, aiming to find a minimal set of configurations that achieves the same coverage as the full set of runs. We begin with the empty list of configurations. At each step of the algorithm, we pick the interaction that (if we also include the coverage of all subsets of that interaction) guarantees the most as-yet-uncovered lines, blocks, etc. Then, we scan through the list to find a configuration that is consistent with our pick. We merge the interaction with the first such configuration we find in the list, or append the interaction to the list as a new configuration if it is inconsistent with all existing configurations. This algorithm will always terminate and cover all lines, blocks, etc., though it is not guaranteed to find the actual minimum set.

Figure 8 summarizes the results of our algorithm. The column labeled 1 shows how many lines, blocks, edges, or conditions are covered by the first configuration in the list. Then

Config #	1	2	3	4	5	6	7	8	9	10
vsftpd										
Line	2,521	18	8	1	1	-	-	-	-	-
Block	2,853	25	9	1	1	-	-	-	-	-
Edge	2,731	50	17	6	1	1	1	-	-	-
Condition	1,132	71	14	9	2	1	1	1	1	-
ngIRCd										
Line	3,148	30	6	6	1	1	1	-	-	-
Block	4,401	50	8	7	4	1	1	-	-	-
Edge	4,390	62	14	8	6	2	2	2	-	-
Condition	1,881	27	23	5	4	1	1	1	1	-
grep										
Line	2,218	171	34	20	5	5	3	2	2	-
Block	2,838	231	46	28	5	5	3	1	-	-
Edge	3,140	366	51	44	18	9	6	6	4	-
Condition	1,810	231	45	25	11	8	7	6	5	1

Figure 8: Additional coverage achieved by each configuration in the minimal covering sets.

column n (for $n > 1$) shows the additional coverage achieved by the n th configuration over configurations 1..($n - 1$). Notice that minimal covering sets range in size from 5 to 10, which is much smaller than the number of possible configurations. This suggests that when we abstract in terms of coverage, in fact the configuration space looks more like a union of disjoint interactions (that can be efficiently packed together) rather than a monolithic cross-product.

We can also see that each subject program follows the same general trend, with most coverage achieved by just the first configuration. The last several configurations very often add only one additional coverage element. This last finding hints that not every interaction offers the same level of coverage; we explore this issue in detail in the next section.

Finally, we also used this algorithm to compute the full *effective configuration space* of each program, which is a set of configurations that ensures that every realizable *path* is executed at least once. The effective configuration spaces of vsftpd, ngIRCd, and grep contained 3,092, 3,518, and 19,301 configurations, respectively. While significantly larger than for the simpler coverage criteria, these numbers are still far smaller than the product of all possible values of configuration options.

5.4 Configuration Space Analysis

Figure 9 visualizes the interactions of each subject program, to help us understand why the minimal covering sets are so small. These graphs show interactions based on line coverage. Because the full set of interactions is too large to display easily, we show only those interactions needed to guarantee 95% of the maximum possible coverage.² In these graphs, a node represents one or more option settings; we merged nodes with common neighbors, listing all settings the node represents. Shaded nodes are standalone interactions; solid edges mark interactions between pairs of nodes; and patterned triangles denote interactions among triples of nodes. The boxes in Figure 9(a) and (b) represent interactions among all options within them. For ngIRCd, the options within this “super node” also interact, collectively, with the other option settings, as indicated. The ngIRCd options are all prefixed with `Conf_.`, and similarly the vsftpd options are prefixed with `tunable_.` We omitted these prefixes

²Diagrams of the full set of interactions can be found in a companion technical report [13].

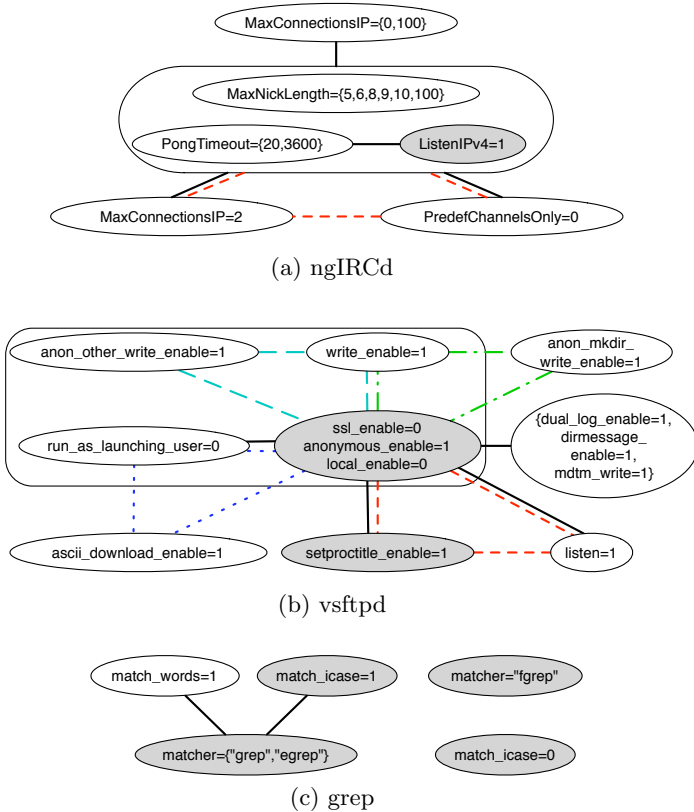


Figure 9: Interactions needed for 95% line coverage. ngIRCd and vsftpd include some approximations.

from the graph to save space.

To unclutter the presentation and to highlight interesting interaction patterns, we made some additional simplifications. For ngIRCd, we merged two values for PongTimeout that had similar but not identical neighbor sets, and similarly for MaxNickLength. For vsftpd, we merged the options in the center node of Figure 9(b) even though they have slightly different neighbors.

The main feature we see in ngIRCd’s graph is the super node in the middle, which contains ngIRCd’s enabling options. We can even see their progression: setting ListenIPv4=1 is the first crucial step that lets ngIRCd accept clients, and it forms a 1-way interaction. Next, setting PongTimeout high enough avoids early termination of client connections, and therefore this option forms a 2-way interaction with ListenIPv4=1. The last enabling option, MaxNickLength, forms a 3-way interaction with the previous two. In the full ngIRCd graph, the full set of these enabling options is similarly connected to most of the nodes in the graph.

Next, considering vsftpd’s graph, we clearly see that all of the interactions involve the enabling options, which appear in the center, shaded node. There are many interactions involving just one additional option setting, such as the three options in the node at the right middle position. These options control the availability of some features, e.g., dirmessage_enable enables the display of certain messages. Moreover, notice that we can combine all the settings in the nodes of Figure 9(b) into one configuration. This helps illustrate why the minimal covering set of configurations for

vsftpd is so small, and why the initial configuration is able to cover so much: one configuration can enable a range of features (writing files, logging, etc.) all at once.

For vsftpd, the full interaction graph is much like the image shown here. The full graph includes a few additional, higher-strength interactions that include the enabling options, plus some low-strength interactions that each guarantee a few additional lines.

Finally, in grep’s graph, notice how few configuration options contributed to 95% of the coverage. These high-coverage interactions of grep have very low interaction strength; there are no interactions with strength higher than two, and four out of the five nodes have 1-way interactions. Also, all values of the matcher option appear in this graph, making this the most important option for grep in terms of coverage. The full configuration space graph of grep contains many more interactions and, interestingly, the important matcher option only takes part in a few interactions in the full graph.

While each program exhibits somewhat different configuration space behavior, we can see that when abstracted in terms of line coverage, many options either do not interact or interact at low strength, and thus we can combine them together into larger configurations. This supports our claim that configuration spaces are considerably smaller than combinatorics might suggest.

5.5 Threats to Validity

Like any empirical study, our observations and conclusions are limited by potential threats to validity. For example, in this work we used 3 subject programs. Each is widely used, but small in comparison to some industrial applications. In order to keep our analyses tractable, we focused on sets of configuration options that we determined to be important. The size of these sets was substantial, but did not include every possible configuration option. The program behaviors we studied included four structural coverage criteria. Other program behaviors such as data flows or fault detection might lead to different results. Our test suites taken together have reasonable, but not complete, coverage. Individually the test cases tend to be focused on specific functionality, rather than combining multiple activities in a single test case. In that sense they are more like a typical regression suite than a customer acceptance suite. We intend to address each of these issues in future work.

6. RELATED WORK

Symbolic Evaluation. In the mid 1970’s, King was one of the first to propose symbolic evaluation as an aid to program testing [9]. Theorem provers at that time, however, were fairly simple, limiting the approach’s practical potential. Recent years have seen remarkable advances in Satisfiability Modulo Theory and SAT solvers, which has enabled symbolic evaluation to scale to more practical problems. Some recent symbolic evaluators include DART [7, 8], EXE [3], and KLEE [2]. There are important technical differences between these systems, e.g., DART uses *concolic execution*, which mixes concrete and symbolic evaluation, and KLEE uses pure symbolic evaluation. However, at a high level, the basic idea is the same: the programmer marks values as symbolic, and the evaluator explores all possible program paths reachable under arbitrary assignments to those symbolic values. As mentioned earlier, Otter is closest in

implementation terms to KLEE.

SE for Configurable Systems. Researchers and practitioners have developed several strategies to cope with the problem of testing configurable systems. One popular approach is combinatorial testing [4, 1], which, given an *interaction strength* t , computes a *covering array*, a small set of configurations such that all possible t -way combinations of option settings appear in at least one configuration. The subject program is then tested under each configuration in the covering array, which will have very few configurations compared to the full configuration space of the program.

Several studies to date suggest that even low interaction strength (2- or 3-way) covering array testing can yield good line coverage while higher strengths may be needed for edge or path coverage or fault detection [1, 5, 10]. However, as far as we are aware, all of these studies have taken a black box approach to understanding covering array performance. Thus it is unclear exactly how well and why covering arrays work. On the one hand, a t -way covering array contains all possible t -way interactions, but not all combinations of options may be needed for a given program or test suite. On the other hand, a t -way covering array must contain many combinations of more than t options, making it difficult to tell whether t -way interactions, or larger ones, are responsible for a given covering array's coverage. Our work attempts to better understand what specific configuration space characteristics control system behavior.

7. CONCLUSIONS AND FUTURE WORK

We have presented an initial experiment using symbolic evaluation to study the interactions among configuration options for three software systems. Keeping existing threats to validity in mind, we drew several conclusions. All of these conclusions are specific to our programs, test suites, and configuration spaces; further work is clearly needed to establish more general trends.

First, we found that we could achieve maximum coverage without executing anything near all the possible configurations. Most coverage was accounted for by lower-strength interactions, across all of line, basic block, edge, and condition coverage. Second, if we packed interactions into configurations greedily, it took only five to ten configurations to achieve this maximal coverage. Third, we also found that in fact it only took one configuration to get the vast majority of the maximum coverage. Finally, by mapping the interactions we gained some insight into why the minimal covering sets are so small. We observed that many options either did not interact or interacted at low strength, and it is often possible to combine different interactions together into a single configuration. Taken together, our results strongly suggest our main hypothesis—that in practical systems, configuration spaces are significantly smaller than combinatorics suggest, and they can be understood as the composition of a small number of interactions.

Based on this work, we plan to pursue several research directions. First, we will extend our studies to better understand how configurability affects software development. Some initial issues we will tackle include increasing the number and types of options and repeating our study on more and larger subject systems. Second, we plan to enhance our symbolic evaluator to improve performance, which should enable larger scale studies. One potential approach is to

use path pruning heuristics to reduce the search space, although we would no longer have complete information. Finally, we will explore potential applications of our approach and results, such as test prioritization, configuration-aware regression testing, and impact analysis.

Acknowledgments

We would like to thank Cristian Cadar and the other authors of KLEE for making an early version of their code available to us. We would also like to thank Mike Hicks, Stephen Magill, and the anonymous reviewers for their helpful comments. The author Charles Song was partially supported by Fraunhofer CESE, USA. This research was supported in part by NSF CCF-0346982, CCF-0430118, CCF-0524036, CCF-0811284 and CCR-0205265.

8. REFERENCES

- [1] R. Brownlie, J. Prowse, and M. S. Phadke. Robust testing of AT&T PMX/StarMAIL using OATS. *AT&T Technical Journal*, 71(3):41–7, 1992.
- [2] C. Cadar, D. Dunbar, and D. R. Engler. KLEE: Unassisted and automatic generation of high-coverage tests for complex systems programs. In *OSDI*, pages 209–224, 2008.
- [3] C. Cadar, V. Ganesh, P. M. Pawlowski, D. L. Dill, and D. R. Engler. EXE: automatically generating inputs of death. In *CCS*, pages 322–335, 2006.
- [4] D. M. Cohen, S. R. Dalal, M. L. Fredman, and G. C. Patton. The AETG system: an approach to testing based on combinatorial design. *TSE*, 23(7):437–44, 1997.
- [5] I. S. Dunietz, W. K. Ehrlich, B. D. Szablak, C. L. M. ws, and A. Iannino. Applying design of experiments to software testing. In *ICSE*, pages 205–215, 1997.
- [6] V. Ganesh and D. L. Dill. A decision procedure for bit-vectors and arrays. In *CAV*, July 2007.
- [7] P. Godefroid, N. Klarlund, and K. Sen. DART: directed automated random testing. In *PLDI*, pages 213–223, 2005.
- [8] P. Godefroid, M. Y. Levin, and D. A. Molnar. Automated whitebox fuzz testing. In *NDSS*. Internet Society, 2008.
- [9] J. C. King. Symbolic execution and program testing. *Commun. ACM*, 19(7):385–394, 1976.
- [10] D. Kuhn and M. Reilly. An investigation of the applicability of design of experiments to software testing. In *NASA Goddard/IEEE Software Engineering Workshop*, pages 91–95, 2002.
- [11] G. C. Necula, S. McPeak, S. P. Rahul, and W. Weimer. CIL: Intermediate language and tools for analysis and transformation of C programs. In *CC*, pages 213–228, 2002.
- [12] A. Porter, C. Yilmaz, A. M. Memon, D. C. Schmidt, and B. Natarajan. Skoll: A process and infrastructure for distributed continuous quality assurance. *TSE*, 33(8):510–525, August, 2007.
- [13] E. Reisner, C. Song, K.-K. Ma, J. S. Foster, and A. Porter. Using symbolic evaluation to understand behavior in configurable software systems. Technical Report CS-TR-4946, Department of Computer Science, University of Maryland, 2009.