Bug Synthesis: Challenging Bug-Finding Tools with Deep Faults

Subhajit Roy Computer Sc. and Engg. IIT Kanpur, India subhajit@cse.iitk.ac.in Awanish Pandey Computer Sc. and Engg. IIT Kanpur, India awpandey@cse.iitk.ac.in Brendan Dolan-Gavitt Computer Sc. and Engg. NYU Tandon, USA brendandg@nyu.edu Yu Hu Computer Sc. and Engg. NYU Tandon, USA yh570@nyu.edu

ABSTRACT

In spite of decades of research in bug detection tools, there is a surprising dearth of ground-truth corpora that can be used to evaluate the efficacy of such tools. Recently, systems such as LAVA and EvilCoder have been proposed to automatically inject bugs into software to quickly generate large bug corpora, but the bugs created so far differ from naturally occurring bugs in a number of ways. In this work, we propose a new automated bug injection system, Apocalypse, that uses formal techniques—symbolic execution, constraint-based program synthesis and model counting—to automatically inject fair (can potentially be discovered by current bug-detection tools), deep (requiring a long sequence of dependencies to be satisfied to fire), uncorrelated (each bug behaving independent of others), reproducible (a trigger input being available) and rare (can be triggered by only a few program inputs) bugs in large software code bases. In our evaluation, we inject bugs into thirty Coreutils programs as well as the TCAS test suite. We find that bugs synthesized by Apocalypse are highly realistic under a variety of metrics, that they do not favor a particular bug-finding strategy (unlike bugs produced by LAVA), and that they are more difficult to find than manually injected bugs, requiring up around 240× more tests to discover with a state-of-the-art symbolic execution tool.

CCS CONCEPTS

Software and its engineering → Software testing and debugging; Formal software verification;

KEYWORDS

Bug Injection, Program Synthesis, Symbolic Execution, Constraint-based Synthesis

ACM Reference Format:

Subhajit Roy, Awanish Pandey, Brendan Dolan-Gavitt, and Yu Hu. 2018. Bug Synthesis: Challenging Bug-Finding Tools with Deep Faults. In Proceedings of the 26th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE '18), November 4–9, 2018, Lake Buena Vista, FL, USA. ACM, New York, NY, USA, 11 pages. https://doi.org/10.1145/3236024.3236084

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. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ESEC/FSE '18, November 4-9, 2018, Lake Buena Vista, FL, USA

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 978-1-4503-5573-5/18/11...\$15.00

https://doi.org/10.1145/3236024.3236084

1 INTRODUCTION

Decades of research has gone into eliminating bugs in software through automated bug-finding tools such as static analyzers, runtime sanitizers, symbolic execution tools, and fuzzers. Despite this ongoing effort, there is a surprising dearth of ground-truth corpora that allows us to evaluate the *efficacy* of such tools: most existing corpora are small, do not come with triggering inputs, or feature bugs that are often unrealistic. Moreover, the value of any individual dataset drops over time as tools adapt to it. As a result, in many cases we are forced to judge bug-finding tools on how many previously unknown bugs they find—which leaves us in the dark about how many they missed.

The lack of ground-truth datasets also means that it is difficult to perform large-scale studies of bug discovery. For example, we cannot run bug-finding tools on corpora of millions of bugs and then attempt to draw conclusions about their relative strengths and weaknesses, or statistically correlate features of bugs and programs with their difficulty of discovery.

In recent work, two new systems, LAVA [14] and EVILCODER [27], have sought to address the need for ground-truth corpora through *automated vulnerability addition*. Briefly, these techniques take existing programs and seed them with vulnerabilities, either by adding new, vulnerable code, or by identifying and removing safety checks to make existing code vulnerable. These systems are an important step forward.

We consider that injected bugs should have the following properties:

- Fair An injected bug is *fair* when it is possible to unearth it by practical bug detection techniques. For example, a bug guarded by a famous mathematical theorem (say, Fermat's Last Theorem) is not fair as it requires the proof of a difficult mathematical theorem to detect the bug; common bugs in programs do not resemble this.
- Deep An injected bug is deep if it requires a long sequence
 of data and control flow conditions to be met for it to trigger.
 A bug guarded by a single branch condition is, for the same
 reason, not a challenging bug.
- Uncorrelated Multiple injected bugs must be uncorrelated; that is, finding one of the bugs by a tool should not increase (or decrease) the chances of catching the other injected bugs.
- Reproducible An injected bug must come with a triggering input that proves the existence of the bug.
- Rare The bug should be triggered on a very small fraction of all possible program inputs.

Considering these properties, we find that existing bug injection techniques can be improved in several ways. EvilCoder, for example, cannot produce triggering inputs, and hence fails to be *reproducible* (related techniques, such as mutation testing [21], also fail to satisfy

this requirement). And we find that bugs injected by LAVA, although rare, uncorrelated, and reproducible, fail to be fair and deep: the triggers used (a comparison against a 32-bit "magic" constant) are unusually difficult for techniques such as random testing to find, and although the bugs manifest deep within programs, the injected guard is a single branch that can be systematically targeted [32].

In this work, we introduce a new technique for bug injection, based on symbolic execution, program synthesis and uniform sampling. We build our ideas into a tool, APOCALYPSE, and use it to introduce bugs in thirty Coreutils programs. Apocalypse uses constraint based program synthesis to embed a transition system, that we refer to as the Error Transition System (ETS), on a judiciously chosen program path. When the program is executed, the ETS is advanced at certain locations along this execution path, leading to a crash if the final state is reached. The state transitions on the ETS are guarded by carefully synthesized predicates that ensure that only a few executions can successfully reach the final state, and, therefore, trigger the bug. We do so by enabling the synthesis engine to perform a multi-variate hill climbing on the space of predicates at the transition locations—searching for predicates that prevent most executions from reaching the bug location. We estimate the set of inputs that a predicate "blocks" from reaching the error location by model counting (approximated by uniform sampling). Model counting and uniform sampling are not mature technologies yet, but are quickly making progress with some significant contributions in recent years.

We use Apocalypse to inject multiple bugs in thirty Coreutils programs, and then attempt to detect these bugs using state-of-theart symbolic execution (KLEE [8]) and greybox fuzzing (AFL [1]) tools. Many of the injected bugs were discovered by the above tools, showing that the injected bugs were fair: KLEE and AFL were able to discover 31% and 38% of the bugs (respectively); at the same time, many bugs were elusive, showing that these bugs can act as subjects for further research: 47% of the bugs could not be discovered by either of the tools. Similar to real bugs, different bugs showed affinity to different tools: out of the 53% bugs discovered, 15% of the bugs could be discovered only by KLEE while 22% of the bugs could be found only by AFL. Bugs synthesized by Apocalypse needed about 240× more tests to be discovered than the manually seeded bugs (on our benchmarks). We also compared our tool with LAVA and found that the bugs injected by LAVA tend to be biased to one of the bug-finding tools: about 80% of the bugs injected by LAVA were discovered by KLEE while only about 41% of these bugs were discovered by AFL; bugs injected by APOCALYPSE, on the other hand, responded almost uniformly to both KLEE and AFL, showing the bug corpora produced by APOCALYPSE do not favor a particular bug-finding strategy.

As most bug-finding tools do not target any attributes of real-world bugs but are driven solely by code-coverage (like KLEE and AFL), Apocalypse synthesizes bugs that challenge the code-coverage heuristics of these tools, stressing their competence at uncovering deep paths.

The contributions of this paper are as follows:

We propose a symbolic execution-based strategy to automatically inject fair, deep, uncorrelated, reproducible and rare bugs in programs;

- We propose a model counting-based strategy to reduce the number of bug inducing inputs to make the injected bugs difficult to find;
- We build our ideas into a tool, APOCALYPSE, and use it to inject bugs into Coreutils programs;
- We attempt to gauge the *quality* of the bugs synthesized by APOCALYPSE using a symbolic execution engine (KLEE) and a greybox fuzzer (AFL).

2 OVERVIEW

Our bug injection system, APOCALYPSE, begins with a concrete input and a program trace induced by that input. This input, which can be taken from the program's test suite or (as in our current implementation) discovered through symbolic execution, serves as a path along which we add one or more bugs to the program. One can envisage bug synthesis as a game between the injector (who wishes to add hard-to-find bugs to the program) and the bug-finder (who would like to find the bugs added by the injector), the concrete input serving as a source of asymmetric advantage in favor of the injector: armed with a concrete input, the injector has knowledge of an entire program path and all dynamic values along that path, whereas the bug-finder must search for the same program path in the space of all program paths.

APOCALYPSE embeds a state machine, which we refer to as the *error transition system* (ETS), within the subset of the program statememts described by the trace; this ETS is designed such that the execution of the program on a pre-selected fault-revealing input would incrementally advance this state machine towards an error state (i.e., the program point where the injected bug will manifest). A transition in the state machine is triggered whenever the current values of pre-existing program variables (that are in scope at the current location) meet certain conditions. To create a bug satisfying the requirements described in §1, these conditions must be *simple* (cause minimum perturbation to the dataflow and control-flow of the existing program), *non-trivial* (are not always true or always false), and *useful* (satisfied by a relatively small number of inputs).

APOCALYPSE achieves these goals by using *program synthesis* to create candidate predicates from variables that are in scope at different points (*transition points*) along the trace. Not all points in the trace are equally promising as transition points: we scan the trace looking for program points that are deep in the program call graph, guarded by many branch conditions, and have many variables in scope. To ensure that the synthesized constraints that trigger each transition meet our requirements, we use *model counting* to estimate the number of solutions to the conjunction of ETS constraints so far, and iteratively improve the constraint set by reducing the number of possible solutions.

The state machine is tracked using the global program variables. We encode the state machine to match the lexical charateristics of the subject program: for example, if the program primarily manipulates integer values, we use entities of the integer type (integer variables, elements of integer arrays, integer fields in structures etc.) to track the state (this is the encoding supported by our current prototype). But we could also track the state using string matching, or the position of some node in an aggregate data structure (such as a list or a tree), depending on what data structures and operations

```
void ALIM ()
othCap = climb + altVal;
L6: /* if (ownRate < otherAlt && state == 6)
           state = 19; */
int InhibitBiasedClimb ()
int up, down;
up = upSep + 100 + altVal;
down = upSep + OtherTrackedAlt;
L16: /* if (othCap < climb && state == 16)
            state = 6; */
return (climb ? up : down);
void main()
input \, (\, curSep \, , \, \, ownAlt \, , \, \, ownRate \, , \, \, otherAlt \,
       altVal, upSep, downSep, othCap, climb);
L0:
   state = 16; */
upPref = InhibitBiasedClimb() + downSep;
upCross = ownAlt + otherAlt;
ownRate = ownRate + curSep;
ALIM();
 if (uppref > 5500) {
   result = altVal;
L19: /* if (climb == result && state == 19)
            state = 21; */
L21: /* if (othAlt < upCros && state == 21)
            state = 30;
upCross = ownAlt - otherAlt;
L30: /* if (state == 30) assert(0); */
```

Listing 1: Program with synthesized bug: the statements commented in *green* are statements inserted by APOCALYPSE (as per the ETS in Figure 1)

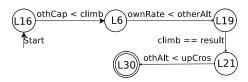


Figure 1: ETS for the program in Listing 1

present in the subject program. We give some concrete examples of possible state machine encodings in §3.7.

Finally, we create a buggy version of the program by adding, at each transition point, a snippet of code that checks one of our ETS transition predicates and then, accordingly, advances the state machine. When the state machine reaches its accepting state, we trigger (a pre-defined) buggy behavior. In our current implementation, we force an assertion failure, simply by adding an assert(false); however, other buggy behaviours (like out of bounds memory accesses, double free, memory leaks, dereferencing of freed pointers, floating-point exceptions) can also be added depending on the class of the bug detector employed.

Listing 1 shows a program with an injected bug; the statements synthesized by Apocalypse are shown commented in *green*; the program is instrumented with the ETS shown in Figure 1.

Algorithm 1 APOCALYPSE

```
1: procedure Main(\mathcal{P})
2: \Delta \leftarrow \text{IdentifyTrace}(\mathcal{P})
3: L \leftarrow \text{IdentifyTransitionPoints}(\Delta)
4: sym\_pc, \Lambda \leftarrow \text{SymbolicExec}(\Delta, L)
5: \Omega \leftarrow \text{SynthesizeETS}(sym\_pc, \Lambda)
6: \mathcal{P}' \leftarrow \text{InstrumentETS}(\mathcal{P}, \Omega)
7: return \mathcal{P}'
8: end procedure
```

3 ALGORITHM

We define a *program trace* (or simply a *trace*) Δ as a sequence of dynamic instructions. We assume each trace Δ to have a triggering input ip, that causes the program to execute the given trace, and a symbolic path constraint, sym_pc, that encodes the conditions on the inputs that would follow the given trace. Given a map $\Gamma: V \mapsto \mathcal{E}$ from program variables $v_i \in V$ to symbolic expressions $s \in \mathcal{E}$, we use the notation $\Gamma[\phi(v_1, \ldots, v_n)]$ to denote the symbolic constraint formed by replacing each $v_i \in V$ by the respective symbolic expression from the map Γ .

3.1 Error Transition System

Apocalypse injects bugs in programs by interweaving an *Error Transition System* along a path in the program. The ETS is a tuple $(L, \mathcal{P}, \delta, l_0, l_{buq})$ where:

- *L*, the set of states, corresponds to program locations that drive a transition on the ETS;
- \mathcal{P} is a set of all predicates that can be constructed using the program variables, array elements and structure fields in the current scope;
- δ: L×P → L is the transition function that dictates the transition on the ETS, given a predicates p ∈ P at a location l_i ∈ L;
- l₀ ∈ L is the initial state of the ETS; the ETS is set to this state at the entry point of the program;
- *l_{bug}* ∈ *L* is the program location that is instrumented by the buggy action (say simulating a program crash).

Figure 1 shows the ETS for the program in Listing 1: the labels L6, L16 etc. (marked in red) show the transition locations where the ETS makes its moves; the program entry point (L0) sets the ETS to the initial state (L16). The transition on location L16 is guarded by the predicate (othCap < climb); if this predicate holds, the ETS transitions to state L6. Finally, if the execution can drive the ETS to the final state (L30), an error is raised (say by simulating a crash or violating an assertion).

3.2 Identify a Program Trace

The algorithm driving Apocalypse is shown in Algorithm 1. Given a program \mathcal{P} , the algorithm starts off by using IdentifyTrace(\mathcal{P}) to identifying a $trace\ \Delta$ on which an ETS will be embedded. IdentifyTrace(\mathcal{P}) uses symbolic exploration to collect multiple possible interprocedural paths in the program and selects a path based on the following parameters:

• Complexity of the path: We prefer program paths that contain a large number of dynamic instructions, pass through

a large number of procedures and hit a large number of branching instructions. As this path represents the *secret* information that the adversary (bug detection tool) will need to discover, a complex path makes the injected bugs more elusive.

• Number of *useful* variables: This refers to the quality and quantity of the variables (including array elements and structure fields) that are used by the participating instructions along this path. The quality of a variable is dictated by distance of the instruction that defines the variable in the *program dependence graph* from the input statements. In essence, it captures the "complexity" of constructing a required value into this variable from the program inputs. We select paths with abundant *good-quality* variables as these variables are eventually used by the *ETS synthesizer* to construct transition predicates.

3.3 Identify Transition Points in the Program

Next, IdentifyTransitionPoints(Δ) attempts to find good program locations on the error trace to embed ETS transitions. A location is selected if it meets the following criteria:

- Abundant "useful" variables are available at that program location;
- The program location is deep in the call graph, making it hard for bug detection tools to reach this location;
- The program location appears deep in the control dependence graph; a location deep in the control dependence graph is guarded by multiple predicates, making reachability challenging for bug-detection tools.

The above metrics on identifying a trace and transition locations can be tuned to inject bugs of varying degrees of difficulty, thereby allowing one to gauge the effectiveness of different bug detection techniques. In this project, we have attempted to inject bugs that are hard to find; we plan to investigate on the above questions in future work.

3.4 Collect Symbolic Constraints

In the next phase, we run a symbolic execution engine on trace Δ to collect the following:

- Symbolic Path Condition (sym_pc): The path condition (sym_pc) for the trace (Δ) contains a symbolic summarization of all possible input values that would drive a program execution along Δ.
- Symbolic Expression Dictionary (Λ): This dictionary
 Λ: L → (V → ε) maps each identified transition location
 l_i ∈ L in the program to a dictionary of symbolic expressions
 ε for each program variable v ∈ V.
- Concrete Value Dictionary (C): This dictionary
 C: L → (V → v) maps each identified transition location
 l_i ∈ L in the program to a dictionary of concrete values v observed for each program variable v ∈ V along the execution trace.

3.5 Synthesize the Error Transition System (ETS)

In this phase, we use constraint solving to synthesize an Error Transition System (ETS) that can be embedded in the program. The synthesis algorithm is shown in Algorithm 2.

Synthesis of the ETS essentially involves identification of the transition predicates that guard the automata transitions. The identified predicates should satisfy the following properties:

- Simple: The predicate should be simple to compute so as to not change the control-flow and dataflow behaviour of the existing program by much.
- Non-Trivial: The predicates should be non-trivial; for example, $(x \ge x)$, (x + 42 > x) etc. should not be produced.
- Useful: The predicate should effectively reduce the number of inputs that could trigger the bug.

For all the transition locations $l_i \in L$, let $pred : L \to \mathsf{Predicate}$ denote a dictionary of the predicates synthesized, such that pred[i] is the predicate at the i^{th} location (l_i) . This map is initialized to $pred : \lambda L$.true.

We synthesize the predicates for the different locations in a round-robin manner; each predicate, pred[i], is synthesized subject to the current values of all other predicates. The predicate for the k^{th} transition location (denoted as $(v_1 \ op \ v_2)^k$) is synthesized using the following **synthesis condition**:

[Equation Synth]

$$(v_1 \ op \ v_2)^k \equiv \exists_{\alpha_1, \dots, \alpha_n} \ (sym_pc \land \prod_{l_i \in L, i \neq k} pred[i]$$

$$\land \llbracket \ C_k(v_1) \ op \ C_k(v_2) \ \rrbracket)$$

$$\land \lnot \llbracket \ \Lambda_k(v_1) \ op \ \Lambda_k(v_2) \ \rrbracket$$

where

$$\llbracket \mathcal{V}_1 \text{ op } \mathcal{V}_2 \rrbracket = \begin{cases} \llbracket \mathcal{V}_1 \rrbracket & < \llbracket \mathcal{V}_2 \rrbracket & \text{for op } =' <' \\ \llbracket \mathcal{V}_1 \rrbracket & \leq \llbracket \mathcal{V}_2 \rrbracket & \text{for op } =' \leq' \\ \llbracket \mathcal{V}_1 \rrbracket & = \llbracket \mathcal{V}_2 \rrbracket & \text{for op } =' =' \end{cases}$$

The above constraint synthesizes a guard $(v_1 \ op \ v_2)$ for the k^{th} transition location if there exists a feasible execution (i.e. feasible values of the input symbolic variables $\alpha_1, \ldots, \alpha_n$) that meets the above condition. Let us explore each term in the synthesis condition:

- The first term ensures that (by satisfying *sym_pc*) the synthesized path preserves the seed execution path as the fault-revealing run, and the respective input is secured as the triggering input;
- The second term $(\prod_{l_i \in L, i \neq k} pred[i])$ ensures that this path satisfies the guard conditions synthesized at all other transition locations thereby creating an *augmented* path that is traversed by fewer inputs;
- The next term ($[\![C_k(v_1) \ op \ C_k(v_2) \]\!]$) synthesizes a guard condition by searching for variables v_1 and v_2 , and an operator op, such that the concrete value of this predicate with the current values of v_1 and v_2 (corresponding to the seed input) is consistent with the seed execution path; this ensures that

Algorithm 2 Synthesis Algorithm

```
1: procedure SynthesizeETS(sym_pc, \Lambda, C, V)
          for all l_i \in L do
2:
3:
              pred[i] \leftarrow true
               sols[i] \leftarrow \{\}
 4:
          end for
5:
          tries \leftarrow 0
6:
7:
          while tries < MAX TRIES do
8:
               tries \leftarrow tries + 1
9:
               for all l_k \in L do
                   \Psi \leftarrow (v_1 \ op \ v_2)^k \wedge \prod_{\vec{\beta} \in sols[k]} \neg \vec{\beta}[(v_1 \ op \ v_2)]
10:
                   if IsSAT Ψ then
11:
12:
                        p \leftarrow \Psi[(v_1 \ op \ v_2)]
                    else
13:
                        continue
14:
                    end if
15:
                    \Phi \leftarrow sym\_pc \land \prod_{l_i \in L, i \neq k} pred[i]
16:
                    count_{old} \leftarrow ModelCount(\Phi \land pred[k])
17:
                    count_{new} \leftarrow ModelCount(\Phi \land p)
18:
                   \textbf{if } count_{new} < count_{old} \textbf{ then } pred[k] \leftarrow p
19:
                    sols[k] \leftarrow sols[k] \cup \Psi[\alpha_1, \ldots, \alpha_n]
20:
21:
               end for
          end while
22:
23: return pred
24: end procedure
```

the injected bug would continue to be triggered by the seed input with this freshly-synthesized guard condition;

• The final term $(\neg \llbracket \Lambda_k(v_1) \text{ op } \Lambda_k(v_2) \rrbracket)$ ensures that there exists a feasible execution along the false path of the synthesized guard; this term is designed to ensure that the symbolic values of the variables v_1 and v_2 are capable of generating an execution along the false branch of the guard condition, thereby preventing generation of trivial predicates (that always evaluate to true).

Now, we lay out the complete synthesis algorithm in Algorithm 2. Lines 3-4 initialize the dictionaries pred (that remembers the current selection of the transition predicates) and sols (that remembers the set of all solutions seen earlier). Then, the procedure enters into an iterative refinement loop to inductively search for good guards for transition predicates: for each transition location $l_k \in L$, Apocalypse tries to find a feasible guard $(v_1 \ op \ v_2)$ (as per Eqn Synth). To ensure monotonicity, it searches for a solution (line 10) while ensuring that any new solution does not include a solution that we have seen earlier (cached in sols[k]).

If Ψ is satisfiable, the predicate is extracted from the model associated with Ψ (lines 11-12); else we move to the next location. The *pred* and *sols* dictionaries are finally updated at lines 19-20 as per the new solution found.

Figure 2 provides a simplified view on the operation of our synthesizer: assuming Trigger as the seed input, the synthesis constraint attempts to search for a point P1 and predicate (denoted by the line) that divides the input space into two partitions: *white*-region that would induce the bug, and *blue*-region that would not. The existence of the point P1 is important to prevent generation of trivial predicates that do not divide the input space (say as lines that are tangents to the input space). Further, in the next iteration,

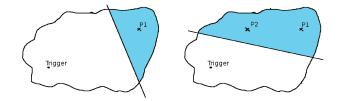


Figure 2: Input space pruning

```
int approxModelCount(oldPred,newPred) {
   /* "sample" is a set of inputs constructed by
   * uniform sampling on the path condition */
   while (inp = sample.next()) {
   if (oldPred[inp]) countOld++;
   if (newPred[inp]) countNew++;
   }
   return (countNew < countOld);
}</pre>
```

Listing 2: Approximate model counting

we further *shrink* the bug-inducing region by searching for another point P2 from the *bug-inducing region* such that the new predicate separates Trigger from both P1 and P2.

To estimate the usefulness of the synthesized predicates, we perform a *hill-climing search* over the multi-variate predicate space corresponding to each location (lines 16-19). This search uses a model-counter to estimate the number of feasible inputs corresponding to the newly synthesized predicate and the older predicate (cached in pred[k]); we select a predicate that maximally *shrinks* the space of bug inducing inputs. The search is designed similar to a Gibbs sampler [9] for multi-variate problems, wherein we make the decision about one variable conditioned on the current values of every other variable.

As invoking a model-counter twice in each iteration is quite expensive, in our implementation, we approximate the relative usefulness of the predicates by performing uniform sampling on sym_pc to create a sampled space of inputs (tests) that follow the same path as the seed input. We, then, generate and execute a program on these tests to count the number of inputs satisfied by the competing guard predicates; we show a sketch of our generated program in Listing 2.

An important design decision was to ensure that our algorithm generates the guard predicates such that the seed input ends up as the trigger for the bug. This is crucial as the symbolic execution engines often *concretize* parts of the execution state (for instance, return values from external library calls, results of floating-point operations, effects of system calls etc.). The map Λ , then, captures only an incomplete symbolic model, potentially leading to path divergence [18].

Let us explain the problem: in Listing 3, assuming sqrt() as an external function, its output will concretized. Hence, the path condition (of an execution where both branches evaluate to true) would be an *incomplete* (a > 42) (rather than (a > 42) \land (c > 10)). If we desire inputs that would trigger the bug by hitting locations L0, L6, and L10, we would need to solve the respective path condition (a > 42) for which a constraint solver can return a solution a = 95, z = 0 (unaware of the branch constraint (c > 10)). This

```
1 L0:
    /* state = 6*/
    make_symbolic(a, z);
    if(a > 42) c = sqrt(z);
    if(c > 10){
        L6:
        /* if(c < a && state == 6) state = 9*/
    }
    L10:
    /* if( state == 9) assert(0)*/</pre>
```

Listing 3: Problem of path divergence

input fails to trigger the bug as the program *diverges* to a different path at the branch "if (c > 10)". We handle this problem by maintaining the seed input (which is consistent with all concretizations) as the trigger for the bug.

For simplicity, the above algorithm assumes that each program location is hit at most once. In our implementation, if a program location is hit multiple times (say in a loop or a procedure), we use only the first few (bounded) instances when the location is reached as possible transition points.

3.6 Embed the Synthesized ETS in the Program

In the final phase, Apocalypse embeds the synthesized ETS in the program by instrumenting the transition locations with guarded state transitions (dictated by the synthesized ETS). Listings 4 and 5 show two possible instrumentation schemes: Listing 5 is a better scheme as it avoids creating path explosion, and hence, creates buggy programs that are closer to the input program in terms of the total number of paths.

Different instrumentation scheme can be adopted to camouflage the ETS transitions: Listing 6 shows a possible camouflage for the ETS transitions for string-processing programs, and Listing 7 for bit-manipulating programs. As our current prototype was meant to study the properties of our injected bugs for automated bug detection systems (and not human subjects), all our experiments were conducted on the instrumentation scheme shown in Listing 5.

3.7 Running Example

To begin with, APOCALYPSE needs to be provided with a *seed input* that drives the program through a path on which we are interested in inducing a fault; a *good* seed path for bug-injection can be discovered by symbolic execution (see §3.2).

Let us work our way through the program in Listing 1: assume that select the seed inputs as (curSep=1258, ownAlt=897, ownRate=174, otherAlt=7253, altVal=1, upSep=629, downSep=5000, otherRAC=0, climb=1). The symbolic and concrete values selected for the inputs are shown in Table 1.

To embed an *Error Transition System* (ETS) along this path, APOCALYPSE also needs a set of *good* program locations to drive the ETS transitions (see §3.3). Our system identifies the lines marked as L16, L6, L19 and L21 as the *transition* locations.

Armed with the seed inputs and the set of transition locations, APOCALYPSE runs symbolic execution along the seed path to collect the symbolic path condition (sym_pc), and the symbolic and concrete expression maps, Λ and C (see Table 2).

Table 1: Symbolic and Concrete inputs for the trace

Variable	Value	Variable	Value	
curSep	$(\alpha_0, 1258)$	ownAlt	$(\alpha_3, 897)$	
ownrate	$(\alpha_4, 174)$	otherAlt	$(\alpha_5, 7253)$	
altVal	$(\alpha_6,1)$	upSep	$(\alpha_7, 629)$	
downSep	$(\alpha_8, 5000)$	othCap	$(\alpha_{10}, 0)$	
climb	$(\alpha_1 1, 1)$			

Table 2: Symbolic and concrete maps

Variable	Loc	16	Loc6		Loc19		Loc21	
	Sym	Conc	Sym	Conc	Sym	Conc	Sym	Conc
othCap	α_{10}	0	$\alpha_6 + \alpha_{11}$	2	$\alpha_6 + \alpha_{11}$	2	$\alpha_6 + \alpha_{11}$	2
ownRate	α_4	174	$\alpha_0 + \alpha_4$	1432	$\alpha_0 + \alpha_4$	1432	$\alpha_0 + \alpha_4$	1432
climb	α_{11}	1	α_{11}	1	α_{11}	1	α_{11}	1
othAlt	α_5	7253	α_5	7253	α_5	7253	α_5	7253
up	$\alpha_6 + \alpha_7$	730						
	+ 100							
down	$\alpha_5 + \alpha_7$	7882						
uppref					$\alpha_6 + \alpha_7$	5730	$\alpha_6 + \alpha_7$	5730
					+ α ₈ +100		$+ \alpha_8 + 100$	
upCros					$\alpha_3 + \alpha_5$	8150	$\alpha_3 + \alpha_5$	8150
result					α_6	1	α_6	1

APOCALYPSE, now, synthesizes an ETS as follows: for the location L16, it finds a predicate (othCap < climb) to move the transition system by a step. It does so by building a **synthesis constraint** that ensures that the predicate is simple, non-trivial (disallowing predicates like $(upPref \ge downSep)$ that are invariants) and useful (discussed next). Similarly, it synthesizes predicates (ownRate < otherAlt), (upPref < upCros)

and (othCap < climb) for locations L6, L19 and L21.

Next, Apocalypse makes more passes over these locations in a search for better predicates. Attempting another synthesis cycle over L19 (and disallowing the previous solution), it synthesizes a new predicate (climb == result). Now, it checks the model count for:

 $sym_pc: (\alpha_{10} \neq -1) \land (\alpha_{11} \neq 0) \land (\alpha_6 + \alpha_8 > 5400) \land (\alpha_5 - \alpha_3 > 0) \land (\alpha_8 \neq 0) \land (\alpha_6 \neq 0)$

 Ψ_1 : $sym_pc \wedge (\alpha_{11} == \alpha_6) \wedge (\alpha_{10} < \alpha_{11}) \wedge (\alpha_0 + \alpha_4 < \alpha_5) \wedge (\alpha_5 < \alpha_3 + \alpha_5)$, and,

 $\Psi_2 : sym_pc \wedge (\alpha_6 + \alpha_7 + \alpha_8 + 100 < \alpha_3 + \alpha_5) \wedge (\alpha_{10} < \alpha_{11}) \wedge (\alpha_0 + \alpha_4 < \alpha_5) \wedge (\alpha_5 < \alpha_3 + \alpha_5).$

In this case, it finds that the model count of Ψ_2 is smaller than that of Ψ_1 , and hence it goes about replacing the older predicate [upPref < upCros] by the newer [climb == result] predicate. On the other hand, if the count of Ψ_1 was smaller, it would have rejected it and persisted with the older predicate. This hill climbing over the multi-variate space of predicates at the different locations allows us to "shrink" the space of inputs that would trigger the bug. Table 3 shows the set of all predicates produced by Apocalypse, with the ones finally selected marked in *blue*. In our experiments, this procedure increased the bug detection time of the injected bugs by about 390× (on AFL). The synthesized ETS is shown in Figure 1.

Finally, the generated ETS is inserted into the existing code. We show the statements injected by Apocalypse as comments (in *green*) in Listing 1; these statements drive the program to a crash at L30.

```
int state:
                        int state:
                                                            char s[50]=
                                                                                           int state = 0xffff0011;
                                                                   hello \0 world \0 for \0 bug 2
                                                                                           void buggy() {
void buggy() {
                        void buggy() {
                                                            void buggyFunction(){
                                                                                             if (p1 && (state
if (p1 && state
                  0)
                         state += 3*(p1 *
                                          !(state - 0));
                                                             if (p1 && !(cmp(s, "hello
                                                                                                  state = 5;
                                                                += (strlen(s)+1);
                          state += 5*(p2 * !(state - 3));
                                                                                                   = 0 \times ff000022;
                                                                                              state
                          state += 3*(p3 * !(state - 8));
if (p2 && state == 5)
                                                             if(p2 && !(cmp(s+l, "world")))
                                                                                             if (p2 && (state
                          state += -12*(p4 * !(state - 11)
state = 9;
                                                                += strlen(s);
                                                                                                  if (p3 && state
                                                             if(p3 && !(cmp(s+l, "for")))
                                                                                                   |= 0 \times ff00 ff00;
                                                                                            if (p3 && (state
                                                                += strlen(s);
if (p4 \&\& state ==-1)
                         if (state == -1)
                                                             if(p4 && !(cmp(s+l, "bug")))
                                                                                                  crash();
                                                              crash();
```

Listing 4: ETS encoding Listing 5: ETS smart encoding for Listing 6: ETS encoding for for integer programs integer programs

string based programs

Listing 7: ETS encoding for bit manipulating programs

Table 3: Synthesized predicates at ETS locations

Loc	Predicates	Loc	Predicates
L16	othCap < climb	L6	ownRate < otherAlt
L19	climb == result	L21	othAlt < upCros
	upPref < upCros		climb < othCap

EXPERIMENTS

APOCALYPSE is based on multiple tools: it uses Clang [2] for instrumentation (for dynamic analysis for selecting good transition locations as well as for embedding the ETS in the program). We modified Crest [7] for running symbolic execution to collect the symbolic path conditions and the expression maps. The ETS synthesizer uses Z3 [11] for constraint solving. We use a modifed version of Boolector [26] to create SAT encodings of SMT constraints, and QuickSampler [15] for uniform sampling on the boolean path conditions.

For the purpose of our experiments, we insert assert(false) statements at our bug injection points. Our experiments were conducted on a 2 GHz Intel Xeon(R) machine with 12 cores and 32 GB RAM. To understand the quality of the bugs injected by Apoca-LYPSE, we attempted to uncover the injected bugs using two popular bug finding techniques:

- **Symbolic Execution**: We use the state-of-the-art symbolic execution engine KLEE [8] to unearth the bugs. KLEE is run with the default search strategy within a timeout of 1 hr.
- Greybox Fuzzing: Coverage-guided fuzzing tools perform executions on randomly mutated inputs, guided by coverage metrics. We use the popular greybox fuzzer AFL [1] for our experiments, running it with default settings and a timeout of 1 hr.

Our experiments attempt to answer the following research ques-

- **RQ1** Are our automatically synthesized bugs fair?
- **RQ2** Is there any correlation between multiple injected bugs?
- **RQ3** Are the bugs injected by Apocalypse reproducible?
- **RQ4** Are our bugs deeper and rarer than manually seeded bugs?
- RQ5 What is the effect of sampling on the difficulty of an injected bug?
- RQ6 How does Apocalypse compare with state-of-the-art bug injection tools?

4.1 RQ1: Fairness of our Synthesized Bugs

We demonstrate that our bugs are fair by employing two state-ofthe-art bug detection tools, KLEE (based on symbolic execution) and AFL (employing greybox fuzzing) for discovering the bugs synthesized by Apocalypse in thirty GNU Coreutils programs [3]. We use Apocalypse to inject four bugs in each program; we then challenge KLEE [8] and AFL [1] discover the bugs (running them for one hour each). The results are shown in Figure 3.

The first bar for each benchmark shows the time spent by AFL to hit each of the bugs, normalized to the time taken to reach the last bug (or timeout when no bug is found). The second bar for each benchmark shows the number of test cases that must be generated by KLEE to find the first test that reveals a bug, normalized to the number of tests required to reach the last bug that it could find.

For example, for the experiments on KLEE (second bars), in the program cat, KLEE is able to find only one bug (so the bar for the first bug reaches all the way to one). In cases where all bugs are found, for example in df, the first bug is found at 15%, the second bug at 32%, and the third at 64% of the total number of testcases needed to find the last bug. Cases where we were unable to find any bug (like test) are shown as timeout (TO).

For the experiments involving AFL (first bars): Because AFL generates tests only for failing executions, we show the amount of time spent to reach a bug (instead of the number of tests). The first bar of each cluster shows how much relative time AFL has invested in finding each bug compared to the last bug. If AFL is not able to find any bug, the whole bar is set to timeout (TO). For example, in the case of cat, the first bug is found at 2%, the second bug is found at 7%, and the third is found at 46% of the time at which the final bug was found.

Overall, KLEE could find 31% of the bugs while AFL found 38.33% of the bugs; 36% of the bugs were found only by one of the two tools while 47% of the bugs could not be found by either. This illustrates common traits exhibited by real bugs:

- Discoverability: State-of-the-art bug detection tools have been successful in dicovering many bugs in large programs. Even for our injected bugs, all in all, 53% of the bugs are discovered by at least one of the tools.
- Elusiveness: Certain bugs are still elusive, showing that these injected bugs (resembling real bugs) can now be employed to stress tools for new bug detection techniques; about 47% of the bugs could not be discovered by either of the tools.

- Affinity to tools: Certain bugs are more likely to be found by one type of technique than by another; 22% of the bugs were only discovered by AFL while 15% of the bugs were only discovered by KLEE.
- Variance in tool effort: Some bugs require more effort to be discovered than others; on the discovered bugs, AFL shows a standard deviation of 583 seconds (on a total running time of 1 hour for each program).

4.2 RQ2: Correlation of Bugs Injected by Apocalypse

Figure 3 shows that the number of tests (using KLEE) and the time taken (by AFL) to discover the different bugs is almost uniformly distributed; also, in many cases, even after discovering a few bugs, the tools fail to unearth the rest of the bugs. This shows that there exists almost no correlation among the different bugs injected by APOCALYPSE. Together, KLEE and AFL are able to catch 53% of all bugs; there are 7 programs (out of 30) where none of the tools is able to catch any bug.

4.3 RQ3: Reproducibility

Because Apocalypse generates triggering inputs for each bug it creates, reproducibility is satisfied by design. Nevertheless, we checked that the generated inputs really did trigger each bug, and found that we could reproduce all the injected bugs.

4.4 RQ4: Comparison with Manually Seeded Bug

To compare with manually seeded bugs, we used the TCAS [13] benchmark. TCAS contains 41 buggy versions, each version containing exactly one manually seeded bug. As the seeded bugs produce an incorrect output (but not a crash or assertion failure), we use KLEE to generate a set of tests; any test that that produces an incorrect output or reaches our injected bug location is designated as a failing test.

We use Apocalypse to inject two additional bugs into each TCAS version. KLEE was able to discover all the injected bugs as well as the manually seeded bugs in all versions except versions 33 and 38. Figure 5 shows the number of test cases KLEE had to generate before hitting the test case that triggers the bug (ploted on the log-scale for the average of three runs of KLEE). This experiment shows the elusiveness of our bugs with respect to the manually seeded ones: on an average, the bugs injected by Apocalypse require 240× more tests than the manually seeded bugs.

Table 4 shows the rarity of our bugs: this table shows the number of generated test cases on which a bug induced a failure. (*F*,*S*) and *MSB* denote the number of failing testcases on our 1st (F) and 2nd (S) injected bug, and manually seeded bug, respectively. On an average, the bugs injected by Apocalypse induce failures on 30× fewer tests over the manually seeded bugs.

4.5 RQ5: Searching for Stronger Transition Predicates

Figure 4 shows how our hill climbing search for guard conditions improves the rarity of the bugs on the different versions of the

Table 4: Number of testcases generated by KLEE that reveal our bugs (F, S) and manually seeded bug (MSB).

V	(F,S)	#MSB	V	(F,S)	#MSB	V	(F,S)	#MSB
1	(251, 1)	232	14	(1, 1)	39	27	(1, 19)	659
2	(1, 1)	160	15	(1, 19)	658	28	(2, 2)	302
3	(1, 1)	41	16	(1, 1)	43	29	(1, 1)	97
4	(4, 3)	374	17	(1, 1)	37	30	(1, 1)	56
5	(1, 19)	654	18	(1, 1)	33	31	(1, 1)	35
6	(1, 1)	30	19	(1, 1)	45	32	(1, 12)	81
7	(1, 1)	41	20	(1, 1)	124	34	(341, 2)	1805
8	(1, 1)	39	21	(9, 1)	162	35	(2, 2)	273
9	(117, 1)	883	22	(1, 1)	172	36	(1, 1)	429
10	(1, 1)	418	23	(1, 1)	424	37	(5, 5)	24
11	(1, 1)	1110	24	(156, 1)	937	39	(1, 1)	209
12	(1, 13)	1783	25	(1, 1)	213	40	(1, 1)	345
13	(1, 1)	310	26	(1, 1)	291	41	(1, 1)	577

TCAS program. We conducted the experiment by comparing the bugs generated when we always picked the first predicate found (red line) versus when the hill-climbing searcher is switched on (blue line). The search for good predicates is effective as it increases the time needed to find the triggering inputs for the bugs by about 390× on average.

4.6 RQ6: Comparison with LAVA

We compare APOCALYPSE with LAVA on nine Coreutils programs (shown in Figure 6). Let us discuss the results on the two tools separately:

- KLEE: KLEE uses constraint solving to discover newer paths. In LAVA, as the bug location is essentially guarded by a magic number, it is quite easy for KLEE to use the constraint solver to "guess" this magic number. Hence, on most of the benchmarks, KLEE is able to discover almost all bugs. On an average, KLEE discovers 80% of the bugs.
- AFL: AFL uses random mutations to discover test cases (guided by coverage information). Hence, AFL finds it hard to guess the magic numbers by random mutations, thereby finding many fewer bugs. On an average, AFL discovers 41% of the bugs.

As can be seen, the bugs injected by LAVA, in general, show affinity towards a certain tool (KLEE). Over a set of injected bugs, each bug may show affinity towards a certain tool, but overall, all bugs injected by a tool should be unbiased. For the bugs injected by APOCALYPSE, though a certain bug may be discovered by a certain tool more easily than the other, overall both tools are almost equally effective (30% of bugs discovered by KLEE, 47% of bugs discovered by AFL on these nine programs) at discovering the bugs injected by our tool. This shows that the bugs synthesized by APOCALYPSE are more "natural" than those injected by LAVA as they do not exhibit "artificial" attributes that can be exploited by bug-finding tools.

5 RELATED WORK

The work most directly related to our current work is LAVA [14] and EvilCoder [27]. As we discuss the relationship of our work to these systems in detail elsewhere in the paper, we omit a complete discussion here, noting only that while our system shares the goals of this prior work, we improve upon the state of the art by

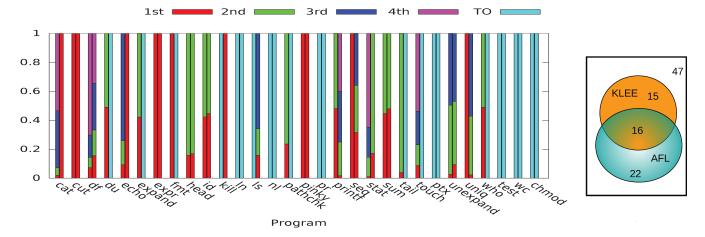


Figure 3: Normalised time (AFL) or testcases (KLEE) needed to find the respective bug; the first bar is for AFL and second for KLEE. The Venn Diagram shows the distribution of the bugs found by KLEE and AFL.

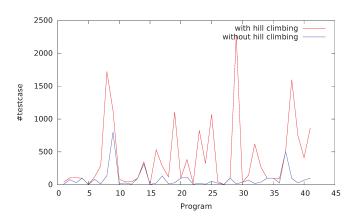


Figure 4: Effect of the hill climbing search on the time to find the bug using AFL in TCAS

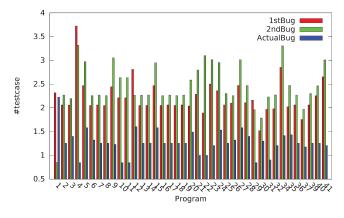


Figure 5: Number of tests generated before hitting a bug (plotted on log-scale for the average of three runs of KLEE)

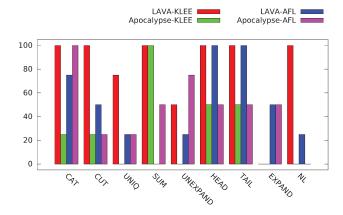


Figure 6: Percentage of bugs found by KLEE and AFL among bugs injected by Apocalypse and LAVA

generating deep bugs that do not favor a particular bug-finding approach.

Bug injection systems are intended to automate the creation of bug corpora; however, it should be noted that there are a number of existing public corpora of buggy programs as well, and studies have been performed to evaluate bug-finders using these systems. For example, Wilander and Kamkar [37, 38] performed a pair of studies using synthetic bugs that evaluated the effectiveness of static and dynamic bug-finding tools. NIST's SAMATE group hosts a collection of buggy software data sets known as the Software Assurance Reference Datasets (SARD) [5]. And in 2016, DARPA hosted an automated bug-finding Cyber Grand Challenge (CGC) [4]; this competition resulted in a collection of 247 programs with known vulnerabilities and triggering inputs, and has been used extensively since its creation for evaluating new bug-finding techniques [29, 31, 32]. The CGC corpus is very high quality, but it is expected that bug-finding software will eventually improve to be able to find all known bugs in the 247 programs. And all of these corpora suffer from one or more of the following issues: they contain few programs or each

```
num input files = argc - optind;
                                                              case 'a': need deadprocs = 1;
                                                              trimmed_name = extract_trimmed_name(utemp_buffer);
case 'r': input_reference = 1;
state12 += (num_input_files > input_reference )
                                                              state12 += (trimmed_name[2] == need_deadprocs)
                * (449- 614) * ¬(state12 -614)
                                                                              * ( 431 - 428 ) * ( ! ( state12 - 428 ));
                      (a) ptx
                                                                                  (b) who
case 'd':
                                                              if (* optarg == '+') from_start = 1;
    output unique = 1
                                                              if (!(xstrtod (optarg, &s, c_strtod) && 0 <= s))
this field = find_field(thisline);
                                                                   sleep_interval = s;
state12 += (thisfield[1] == output_unique) *
                                                              state12 += (*sleep_interval == from_start ) *
                ( 458 - 457 ) * (! ( state12 - 457 ));
                                                                             (621 - 668) * (! (state12 - 668));
                                                                                 (d) tail
                      (c) uniq
case f:
         interactive = 1;
                                                                 (* optarg == '+') from_start = 1;
ok = symbolic link ? symlink (source, dest)
                                                                  (!(xstrtod (optarg, &s, c_strtod) && 0 <= s))
             : link (source, dest)
                                                                   sleep_interval = s;
                                                              state12 += (*sleep_interval == from_start ) *
state12 += ( remove_existing_files < ok ) *</pre>
        ( 136 - 112 ) * (! ( state12 - 112 ));
                                                                             ( 621 - 668 ) * (! ( state12 - 668 ));
                          (e) unia
                                                                                     (f) tail
```

Figure 7: Some of the predicates synthesized by APOCALYPSE for bugs in injected on some Coreutils programs

individual program is small, the bugs may be shallow or unrealistic, or the bugs may not come with triggering inputs.

Another technique for creating bugs is found in the field of *mutation testing* [12, 20–22], in which random *mutation operators* are applied to a program. The resulting (presumably incorrect) program is then run against its test suite in an attempt to judge the robustness of the test suite. In some sense, bug injection is an extension of mutation testing, in that it automatically creates buggy versions of a program. However, the effects of the mutants created by mutation testing are difficult to predict, and do not come with triggering test cases—in other words, they fail to be reproducible (under the definition given in §1). While mutation testing is good for evaluating the quality of a test suite, it is less clear how to apply it to the task of evaluating effectiveness of a bug-finding system such as KLEE [8].

Finally, our bug synthesis strategy is based on the formal techniques like *program synthesis* and *model counting*. Techniques for automatically generating programs have a long history (dating back perhaps as early as 1957, if one includes Church's discussion of the problem of circuit synthesis [10]), but have recently seen a flurry of activity due to the emergence of fast SAT and SMT solvers combined with the work of Solar-Lezama [30], which showed that program synthesis could be cast in terms of satisfiability. Since then, program synthesis has been applied to a wide variety of problems, including automating string processing in spreadsheets [19], heapmanipulations [16, 28] and automated program repair [6, 23, 25, 35]. Model counting [33, 34, 36] and uniform sampling [15, 24] have elicited huge interest due to their applications in bayesian inference and probabilistic programming. Model counting has been successfully employed for probabilistic symbolic execution [17] that assigns probabilities to program paths to aid understanding. We use model counting in a similar context to synthesize a low probability path for the buggy executions.

6 DISCUSSION

Figure 7 shows some the predicates synthesized for a few of the bugs injected by Apocalypse in the Coreutils programs. One can see that the predicates are non-trivial-spanning array accesses, pointer deferences and access to fields of aggregate structures. The asymmetric advantage enjoyed by Apocalypse (as explained in §2) allows it to expend all its reasoning ability on one trace to create challenging bugs, while the bug-finder is required to reason on the whole program. Hence, even with similar compute resources and reasoning ability as the bug-finding tools, Apocalypse is capable of constructing challenging bugs for even state-of-the-art bug-finding tools. We believe that large corpus of injected bugs from Apoca-LYPSE will help us better understand the pecularities and relative merits of different bug-detection tools. There exist threats to validity to our experimental results, in particular from the choice of the subject programs and the seed inputs. We were careful to select a large number of programs and inject multiple bugs; nevertheless, more extensive experiments can be performed.

Acknowledgements. The first author wishes to thank the third author for hosting him at the NYU Tandon School of Engineering, where this work was conceived. Funding for this research was provided under NSF Award #1657199 and ONR Award N00014-15-1-2180.

REFERENCES

- American Fuzzy Lop (AFL) Fuzzer. http://lcamtuf.coredump.cx/afl/technical_details.txt. (visited on 21th January 2018).
- [2] Clang: a C language family frontend for LLVM. http://clang.llvm.org/. (visited on 21th September 2017).
- [3] Coreutils GNU core utilities. https://www.gnu.org/software/coreutils/coreutils. html. (visited on 15th September 2017).
- [4] DARPA. Cyber Grand Challenge. http://archive.darpa.mil/cybergrandchallenge/. (visited on 9th March 2018).
- [5] NIST Software Assurance Reference Dataset Project. https://samate.nist.gov/ SRD/. (visited on 1st January 2018).

- [6] Rohan Bavishi, Awanish Pandey, and Subhajit Roy. 2016. To Be Precise: Regression Aware Debugging. In Proceedings of the 2016 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA 2016). ACM, New York, NY, USA, 897–915. https://doi.org/10.1145/2983990.2984014
- [7] J. Burnim and K. Sen. 2008. Heuristics for Scalable Dynamic Test Generation. In Proceedings of the 2008 23rd IEEE/ACM International Conference on Automated Software Engineering (ASE '08). IEEE Computer Society, Washington, DC, USA, 443–446. https://doi.org/10.1109/ASE.2008.69
- [8] Cristian Cadar, Daniel Dunbar, and Dawson Engler. 2008. KLEE: Unassisted and Automatic Generation of High-coverage Tests for Complex Systems Programs. In Proceedings of the 8th USENIX Conference on Operating Systems Design and Implementation (OSDI'08). USENIX Association, Berkeley, CA, USA, 209–224.
- [9] George Casella and Edward I. George. 1992. Explaining the Gibbs Sampler. The American Statistician 46, 3 (1992), 167–174. http://www.jstor.org/stable/2685208
- [10] Alonzo Church. 1957. Application of Recursive Arithmetic to the Problem of Circuit Synthesis. Summaries of the Summer Institute of Symbolic Logic.
- [11] Leonardo De Moura and Nikolaj Bjørner. 2008. Z3: An Efficient SMT Solver. In Proceedings of the Theory and Practice of Software, 14th International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS'08/ETAPS'08). Springer-Verlag, Berlin, Heidelberg, 337–340. http://dl.acm.org/citation.cfm?id=1792734.1792766
- [12] R. A. DeMillo, R. J. Lipton, and F. G. Sayward. 1978. Hints on Test Data Selection: Help for the Practicing Programmer. Computer 11, 4 (April 1978), 34–41. https://doi.org/10.1109/C-M.1978.218136
- [13] Hyunsook Do, Sebastian G. Elbaum, and Gregg Rothermel. 2005. Supporting Controlled Experimentation with Testing Techniques: An Infrastructure and its Potential Impact. Empirical Software Engineering: An International Journal 10, 4 (2005), 405–435.
- [14] B. Dolan-Gavitt, P. Hulin, E. Kirda, T. Leek, A. Mambretti, W. Robertson, F. Ulrich, and R. Whelan. 2016. LAVA: Large-Scale Automated Vulnerability Addition. In 2016 IEEE Symposium on Security and Privacy (SP). 110–121. https://doi.org/10.1109/SP.2016.15
- [15] Rafael Dutra, Kevin Laeufer, Jonathan Bachrach, and Koushik Sen. 2018. Efficient Sampling of SAT Solutions for Testing. In Proceedings of the 40th International Conference on Software Engineering (ICSE '18). ACM, New York, NY, USA, 549–559. https://doi.org/10.1145/3180155.3180248
- [16] Anshul Garg and Subhajit Roy. 2015. Synthesizing Heap Manipulations via Integer Linear Programming. In Static Analysis, Sandrine Blazy and Thomas Jensen (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 109–127.
- [17] Jaco Geldenhuys, Matthew B. Dwyer, and Willem Visser. 2012. Probabilistic Symbolic Execution. In Proceedings of the 2012 International Symposium on Software Testing and Analysis (ISSTA 2012). ACM, New York, NY, USA, 166–176. https://doi.org/10.1145/2338965.2336773
- [18] Patrice Godefroid, Michael Y. Levin, and David Molnar. 2012. SAGE: Whitebox Fuzzing for Security Testing. Queue 10, 1, Article 20 (Jan. 2012), 8 pages. https://doi.org/10.1145/2090147.2094081
- [19] Sumit Gulwani. 2011. Automating String Processing in Spreadsheets Using Input-output Examples. In Proceedings of the 38th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '11). ACM, New York, NY, USA, 317–330. https://doi.org/10.1145/1926385.1926423
- [20] R. G. Hamlet. 1977. Testing Programs with the Aid of a Compiler. IEEE Trans. Softw. Eng. 3, 4 (July 1977), 279–290. https://doi.org/10.1109/TSE.1977.231145
- [21] Y. Jia and M. Harman. 2011. An Analysis and Survey of the Development of Mutation Testing. IEEE Transactions on Software Engineering 37, 5 (Sept 2011), 649–678. https://doi.org/10.1109/TSE.2010.62

- [22] Richard Lipton. 1971. Fault Diagnosis of Computer Programs. (1971). Student report, Carnegie Mellon Univ.
- [23] Fan Long and Martin Rinard. 2015. Staged Program Repair with Condition Synthesis. In Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering (ESEC/FSE 2015). ACM, New York, NY, USA, 166–178. https://doi. org/10.1145/2786805.2786811
- [24] Kuldeep S Meel, Moshe Y Vardi, Supratik Chakraborty, Daniel J Fremont, Sanjit A Seshia, Dror Fried, Alexander Ivrii, and Sharad Malik. 2016. Constrained Sampling and Counting: Universal Hashing Meets SAT Solving.
- [25] Hoang Duong Thien Nguyen, Dawei Qi, Abhik Roychoudhury, and Satish Chandra. 2013. SemFix: Program Repair via Semantic Analysis. In Proceedings of the 2013 International Conference on Software Engineering (ICSE '13). IEEE Press, Piscataway, NJ, USA, 772–781. http://dl.acm.org/citation.cfm?id=2486788.2486890
- [26] Aina Niemetz, Mathias Preiner, and Armin Biere. 2014 (published 2015). Boolector 2.0 system description. Journal on Satisfiability, Boolean Modeling and Computation 9 (2014 (published 2015)), 53–58.
- [27] Jannik Pewny and Thorsten Holz. 2016. EvilCoder: Automated Bug Insertion. In Proceedings of the 32Nd Annual Conference on Computer Security Applications (ACSAC '16). ACM, New York, NY, USA, 214–225. https://doi.org/10.1145/2991079. 2991103
- [28] Subhajit Roy. 2013. From Concrete Examples to Heap Manipulating Programs. In Static Analysis, Francesco Logozzo and Manuel Fähndrich (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 126–149.
- [29] Yan Shoshitaishvili, Michael Weissbacher, Lukas Dresel, Christopher Salls, Ruoyu Wang, Christopher Kruegel, and Giovanni Vigna. 2017. Rise of the HaCRS: Augmenting Autonomous Cyber Reasoning Systems with Human Assistance. In Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security (CCS '17). ACM, New York, NY, USA, 347–362. https://doi.org/10.1145/3133956.3134105
- [30] Armando Solar-Lezama. 2008. Program Synthesis by Sketching. Ph.D. Dissertation. University of California, Berkeley.
- [31] Nick Stephens, John Grosen, Christopher Salls, Andrew Dutcher, Ruoyu Wang, Jacopo Corbetta, Yan Shoshitaishvili, Christopher Kruegel, and Giovanni Vigna. 2016. Driller: Augmenting Fuzzing Through Selective Symbolic Execution.. In Network and Distributed Systems Symposium (NDSS).
- [32] László Szekeres. 2017. Memory corruption mitigation via hardening and testing. Ph.D. Dissertation. Stony Brook University.
- [33] Marc Thurley. 2006. sharpSAT Counting Models with Advanced Component Caching and Implicit BCP. In *Theory and Applications of Satisfiability Testing -*SAT 2006, Armin Biere and Carla P. Gomes (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 424–429.
- [34] Leslie G Valiant. 1979. The Complexity of Enumeration and Reliability Problems. SIAM J. Comput. 8, 3 (1979), 410–421.
- [35] Sahil Verma and Subhajit Roy. 2017. Synergistic Debug-repair of Heap Manipulations. In Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering (ESEC/FSE 2017). ACM, New York, NY, USA, 163–173. https://doi.org/10.1145/3106237.3106263
- [36] Wei Wei and Bart Selman. 2005. A New Approach to Model Counting. In International Conference on Theory and Applications of Satisfiability Testing. Springer, 324–339.
- [37] John Wilander and Mariam Kamkar. 2002. A Comparison of Publicly Available Tools for Static Intrusion Prevention. In Proceedings of the 7th Nordic Workshop on Secure IT Systems.
- [38] John Wilander and Mariam Kamkar. 2003. A Comparison of Publicly Available Tools for Dynamic Buffer Overflow Prevention. In Proceedings of the 10th Network and Distributed System Security Symposium (NDSS).