

Comparing Application Semantics with S²E

Determining whether busybox and coreutils are actually the same

Mariana D’Angelo

Department of Electrical and Computer Engineering
University of Toronto
Toronto, Canada
mariana.dangelo@utoronto.ca

Dhaval Miyani

Department of Electrical and Computer Engineering
University of Toronto
Toronto, Canada
dhaval.miyani@utoronto.ca

I. INTRODUCTION

While the key idea behind symbolic execution arose almost forty years ago [3], it has only recently become practical due to advances in constraint satisfiability [4] and scalable approaches with mixed symbolic and concrete execution [5], [16]. A 2011 study has shown that symbolic execution tools are beginning to be used increasingly in industrial practice at corporations such as Microsoft, NASA, IBM, and Fujitsu [2]. This proliferation shows the importance of new symbolic execution tools and the impact they will have on the industry as they improve.

S²E is a symbolic execution platform that can operate directly on application binaries for analyzing the properties and behaviour of software systems. It analyzes programs *in-vivo* (the whole environment) within a real software stack (user program, libraries, kernel, drivers, etc.) rather than using abstract models of these layers. It is commonly used for performance profiling, reverse engineering of proprietary software, and finding bugs in user-mode and kernel-mode binaries [1]. S²E is a complex system due to the rich features it has (particularly the fact that it runs inside of a VMM). This complexity, however, comes with the price of being difficult to start using on existing software.

The goal of this experiment is to understand how a complex system such as S²E can be applied to existing software and document the process. In particular, we are interested in seeing whether program inputs and code need to be sanitized in some way to become amenable to dynamic analysis via symbolic execution. Taking inspiration from KLEE [6], we will be comparing the semantics of two implementations of (supposedly) the exact same suite of programs: the GNU Coreutils utility suite and the Busybox embedded system suite.

Comparing the semantics of two implementations can have many real world applications. In the field of computer security, for example, the ability to compare the semantics of a software compiled from source with a pre-compiled binary could reveal malicious or accidental backdoors. In particular, the fact that S²E runs in a VMM would allow for it to be used to test kernel level programs such as file systems (or compare read and write operations different file systems).

We have encountered several challenges in using the S²E platform; one in particular forced us to change the initial goal of our experiment (evaluating two consistency models from S²E) due to the fact that not all of the six consistency models

described in [1] are actually implemented in the platform. Our biggest challenge with this project was trying to modify the programs of interest in such a way that they could be directly comparable. The different approaches we tried will be thoroughly discussed in the Experiments section of this paper along with our own explanations about why they failed or succeeded and what impact these choices had on the program analysis as a whole.

This paper is organized as follows: Section II presents some background regarding symbolic execution, Section III discusses our approach to the experiments, Section IV describes in detail how we attempted to implement our approach, Section V discusses the results of our successful experiment with echo, Section VI presents the status of our implementation, Section VII concludes with some lessons learned, and finally in Section VIII we discuss future work.

II. BACKGROUND

There is much work in symbolic execution for testing [6], [8], [9], [10], [11], [12], [13], [14], [15], [16], [17], [7]. A typical approach for testing is fuzzing [17], which involves generating random inputs for a program in an attempt to exercise all paths in a program (and hopefully hit any buggy code). This is known as concrete execution which involves running the program with deterministic values. Fuzz testing has several limitations, in particular coverage of paths in a program. For example, Listing 1 shows a if statement which is only taken when x is 10. There are 2³² possible values which x could take when an input is randomly generated and as such the probability of actually taking this path is 2⁻³². This example demonstrates that the likelihood of some paths being taken when fuzzing is extremely low, which is the reason for the poor test coverage of concrete execution.

```
1 if (x == 10) {  
2     // path 1  
3 } else {  
4     // path 2  
5 }
```

Listing 1. Code example where fuzzing generally has poor coverage

Symbolic execution, on the other hand, runs the application with symbolic inputs [6] which are initially unconstrained by design. The program executes with these symbolic values and

replaces concrete operations with ones that can manipulate symbolic values. The symbolic execution engine “follows” both paths whenever it encounters a branch on a symbolic value and adds the path condition to a set of constraints known as the “path constraint”. When a bug is encountered in the code a test case (i.e., concrete inputs) can be generated from the path constraint. There are two main limitations of symbolic execution: (1) “the path explosion problem” [1], which is due to the exponential growth of paths in a program caused by conditional branches, and (2) “the environment problem” [6], which is due to interactions with the surrounding environment (e.g, operating system, network, etc.).

In order to deal with the environment problem, one can use a mixture of symbolic and concrete execution (also known as concolic execution [7]). Concolic execution gathers constraints using symbolic execution, but then generates concrete values using a constraint solver in order to address the environment problem. Once these concrete values are generated the program can interact with its environment in a normal fashion without the overheads and problems associated with symbolically executing the environment.

S²E is a symbolic execution engine which can perform concrete execution, symbolic execution, and concolic execution depending on the consistency model chosen. There are six different consistency models in S²E, described below. Figure 1 provides an overview of the models showing how they transition from highly strict to highly relaxed models with some details about what consistency requirements are relaxed for each model. A model is consistent if there exist a globally feasible path through the system for every path explored in the unit. A unit is the block one wishes to analyze and the environment is the rest of the system.

- *Strictly Consistent Concrete Execution (SC-CE)*: No symbolic execution in unit or environment.
- *Strictly Consistent Unit-level Execution (SC-UE)*: Unit is symbolically executed while the environment is executed concretely.
- *Strictly Consistent System-level Execution (SC-SE)*: Unit and environment are executed symbolically - this is the only model that executes the environment symbolically.
- *Local Consistency (LC)*: Similar to SC-UE, but it adheres to constraints that the environment/unit API contracts impose on return values.
- *Relaxed Consistency Overapproximate Consistency (RC-OC)*: Similar to LC, but it ignores the constraints from the environment/unit API contracts.
- *Relaxed Consistency CFG Consistency (RC-CC)*: Similar to SC-UE, but like static analysis it can explore any path in the unit’s inter-procedural control flow graph (even infeasible ones).

There are several differences between S²E and other similar symbolic execution engines. For instance, KLEE uses file system models [6] to avoid symbolically executing the actual filesystem. S²E does not take this approach as writing models is a labour intensive and error-prone process. Cute [7] executes the environment concretely (i.e., without modelling) with a

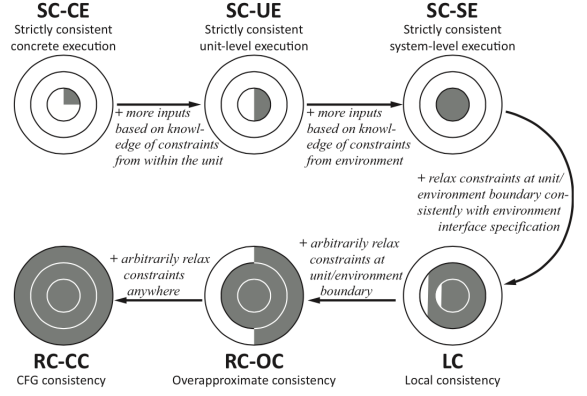


Fig. 1. S²E Consistency Models [1]

consistency model similar to S²E’s SC-SE, but it is limited to code-based selection and one consistency model. S²E does not use compositional symbolic execution [15], a performance optimization which saves results for parts of the program (e.g., a function) and reuses them when that part is called in a different context. With concolic execution everything runs concretely for full systems, however, when the execution crosses program boundaries it may result in lost paths [1]. Due to this, KLEE and CUTE cannot track path conditions in the environment and hence are unable to re-execute calls to enable overconstrained but feasible paths (e.g., malloc does not execute deterministically). DART [16], CUTE [7] and EXE [11] use mix-mode execution (concretely executing some parts and symbolically executing others) to increase efficiency, however, they do not use automatic symbolic-concrete bidirectional data conversion (i.e., automatic conversion between concrete and symbolic values at program boundaries such as the unit and the environment) unlike S²E, which is key to S²E’s scalability and low programmer effort.

Static analysis is performed without executing the program in question and has known limitations such as a high false positive and negative rate (due to infeasible paths) [17] and a lack of runtime environment analysis (as only the application source code is parsed). To reduce the number of false positives tools such as Saturn [19] and bddbdb [20] use a path-sensitive analysis engine. Saturn aims to detect logic programming language bugs by summarizing functions. bddbdb finds buggy patterns in a database where programs are stored as relations. As can be inferred, these tools are language specific and require different implementations for different programming languages. Additionally, using these tools requires learning a new programming language. As a dynamic analysis framework, S²E addresses the limitations of static analysis tools. For example, they directly operate on and analyze binaries whereas static analysis would require disassembly and decompilation. This could involve converting an x86 binary to the LLVM format and running it through an engine like KLEE. Disassembly and decompilation [18] are classically undecidable problems (e.g., disambiguating code from data) [1].

III. APPROACH

The high level approach we took to use S²E for equivalence testing was as follows:

- 1) Run each program with symbolic input (using `s2e_make_symbolic()` from a test application invoking the functions or directly on the binary with option `--sym-args`).
- 2) Compare the output of the Busybox and GNU Coreutils implementation of the same program using `s2e_assert`, which performs a regular assert but also does some cleaning up for S²E and generates sample input which violates the assertion.
- 3) Using the sample input generated, manually analyze the two implementations (possibly using GNU Debugger) to identify semantic differences between the two implementations.

IV. IMPLEMENTATION

In this section we will describe all of the ways in which we tried to symbolically execute the two binaries (from GNU Coreutils and Busybox) to perform equivalence testing.

A. Running the binary directly

Initially, we attempted to run the analysis directly on the binaries for the two implementations of `echo` - having built the two utility suites precisely as suggested by the developers. These results, however, were not easily comparable and did not provide information regarding the semantic differences between the two implementations. From this attempt, it was apparent that running S²E directly on a binary is better suited for simply testing an application with all possible inputs and high coverage, rather than determining semantic equivalence.

Listing 2 shows the command used to execute the binary of GNU Coreutils' `echo`. In order to symbolically execute the binary we must pre-load the `init_env` library (line 1) which intercepts the programs entry point invocation, parses the command line program arguments and configures symbolic execution before invoking `echo`'s main method. We also have to specify some options for our symbolic execution (line 2); `--select-process-code` enables forking the only in the code portion of the current binary, and `--sym-args 0 2 4` generates 0-2 symbolic arguments of length 4 for `echo`. Finally, we use `s2ecmd kill` (line 3) to kill the execution path once `echo` returns to prevent S²E from running forever.

```

1 LD_PRELOAD=/path/to/guest/init_env/init_env.so \
2 /bin/echo --select-process-code --sym-args 0 \
3 2 4 ; /path/to/guest/s2ecmd/s2ecmd kill 0 \
4 "echo done"
```

Listing 2. Command for symbolically executing the GNU Coreutils `echo` binary with S²E

B. Executing the binary from within a tester program

After determining that we would need to invoke both of the implementations from within the same program to compare their semantics, we elected to use the direct approach shown in Listing 3. Here we created a symbolic variable `x` (lines 1 and 2), created the command to invoke `echo` using

`snprintf` (line 3), then enabled S²E's forking (line 5) before using `popen` to execute the command (line 7). We could not use `execv` here because we would not have been able to read what `echo` wrote to `stdout` since `execv` forks a new process, whereas `popen` creates a pipe between the program and the command, and returns a pointer to a stream that can be read from. Later in the code (line 10), we used `fread` to retrieve that stream and read its contents into a buffer. Unfortunately, this technique did not work because S²E was symbolically executing `popen` and so the analysis took too much time, consumed too much memory, and did not yield useful results.

```

1 FILE* pipe; char command[70]; char x[2];
2 s2e_make_symbolic(&x, sizeof(x), "x"); // make x
   symbolic
3 int len = snprintf(command, sizeof(command),
   "echo %s", x);
4 if (len <= sizeof(command)) {
5     s2e_enable_forking();
6     pipe = popen(command, "r");
7     s2e_disable_forking();
8 }
9 ...
10 char buf[100];
11 fread(buf, 1, 1024, pipe);
```

Listing 3. Code to symbolically invoke the GNU Coreutils `echo` binary from a tester program

C. Creating libraries for each implementation

[DHAVAL] Tried to make the `echo` implementations into shared libraries so they could be called directly from our test program (shared library vs. static library stuff) - make sure you discuss why we needed a shared/static library in the first place

D. Redirecting stdout to a buffer

Once we had successfully managed to invoke the two implementations of `echo` from our tester program we needed to retrieve the output of the programs in order to compare their semantics. As both implementations write output directly to the shell we elected to redirect `stdout` using the commands in Listing 4. Although this worked for GNU Coreutils, the invocation of Busybox's implementation of `echo` caused our buffer to be corrupted. Upon investigation, we discovered that although GNU Coreutils uses `stdio`, Busybox avoids using it as the developers do not assume that `stdout` is actually open (and failed writes fill up the input buffer which cannot always be flushed), so they use `writew` to output the result.

```

1 char buf[128];
2 freopen("coreutils.out", "w", stdout);
3 setbuf(stdout, buf);
```

Listing 4. Code snippet for redirecting `stdout` to a buffer

At this point we decided to simply write the result to a file (omitting lines 1 and 3 of Listing 4), which solved our corruption issues. However, when we tried to symbolically execute the program we discovered that S²E was only executing

concretely! Upon further inspection it became apparent that S²E concretized the value of the buffer because our symbolic program argument was only being read from and not written to. Encountering this roadblock, we elected to look into KLEE’s methodology for their comparison of the two utility suites. We discovered that KLEE was symbolically modelling stdout (command line option `-sym-stdout`). Since S²E is built upon KLEE we tried to add this flag into the `kleeArgs` portion of the configuration file used by S²E, however this option did not exist, most likely due to the fact that S²E is intended as an in vivo platform that uses a real software stack (i.e., it avoid modelling parts of the system).

E. Modifying the implementations to not use stdout

[DHAVAL] Final approach: modifying coreutils echo and busybox echo to return a string instead of printing to stdout [SUCCESS!]

V. EXPERIMENTS

For our experiment to check the semantic equivalence of two implementations of `echo` we used GNU Coreutils version 8.23 and Busybox 1.22.1 (the latest versions at the time of this writing).

Once we made changes to both implementations of `echo` to return strings rather than output the result and created shared/static libraries to invoke these functions from our equivalence checking program, we were able to run both programs symbolically and compare their results using the code in Listing 5.

```
1 [INSERT CODE HERE]
```

Listing 5. Program used to perform equivalence testing on GNU Coreutils and Busybox

Our analysis produced the results similar to those shown in listing 6 which provides the constraints for that state (line XX) and also a sample test case that satisfies those constraints (line XX). Such output exists for all states including those that exit normally (i.e., do not fail the assert). There were five assertion failures in our output falling under two categories: those that ended with `0x00` (i.e., the `NULL` character), and those that ended with extended ASCII values like `^` (expressed as a signed integer `-0x65`).

```
1 [INSERT CODE HERE]
```

Listing 6. Partial output of equivalence testing for GNU Coreutils and Busybox

Using the sample input provided by S²E we manually analyzed the two implementations of `echo` to see why their were differences. We determined that although both implementations of `echo` print a space between strings found in `argv`, the conditions leading to printing a space are different. In GNU Coreutils’ implementation, if there is at least one string to echo (check `argv > 0`) a space is added after the current word. Conversely, in Busybox’s

implementation a space is only printed after a string if there is another string to print after the current one (check `***argv != NULL`).

We also had a personal experience of S²E finding a bug due to an implementation error when changing the implementation of GNU Coreutils’ `echo` to return a string rather than output the result. Here we had accidentally changed the code which prints a newline after the echo is completed with code to print a space. Luckily, S²E was able to bring this error to our attention and we fixed this before running our actual experiments.

VI. STATUS

[DHAVAL] Current status of implementation.

- After several failed attempts we decided to modify the utilities and got S2E working with both implementations of echo.

VII. CONCLUSION

[DHAVAL] Did your results meet expectations?

- Summarize intro (what did we do, why did we do it)
- Lessons learned: lack of documentation and an incomplete understanding of how S2E symbolically executes led to a lot of mistakes.

VIII. FUTURE WORK

[DHAVAL]

- Perform these tests on the rest of the utilities in busybox and gnu coreutils.
- Try testing other programs that don’t use stdout and determine if they present any different challenges.
- Figure out how to change the VMM parameters for S2E (e.g., give it more memory). Discuss how even building from source had its problems as we could not transfer files from host to guest when we built S2E manually.

REFERENCES

- [1] CHIPOUNOV, V., GEORGESCU, V., ZAMFIR, C., AND CANDEA, G. Selective symbolic execution. In Workshop on Hot Topics in Dependable Systems, 2009.
- [2] CADAR, C., GODEFROID, P., KHURSHID, S., PASAREANU, C., SEN, K., TILLMANN, N., AND VISSER, W. Symbolic execution for software testing in practice preliminary assessment. In ICSE Impact11, May 2011.
- [3] BOYER, R. S., ELSPAS, B., AND LEVITT, K. N. SELECT - a formal system for testing and debugging programs by symbolic execution. SIGPLAN Not., 10:234245, 1975.
- [4] DE MOURA, L. AND BJERNER, N. Satisfiability modulo theories: introduction and applications. Commun. ACM, 54:6977, Sept. 2011.
- [5] CADAR, C. AND ENGLER, D. Execution generated test cases: How to make systems code crash itself. In SPIN05, Aug 2005.
- [6] CADAR, C., DUNBAR, AND ENGLER, D.R. KLEE: Unassisted and automatic generation of high-coverage tests for complex systems programs. In Symp. on Operating Systems Design and Implementation, 2008.

- [7] SEN, K., MARINOV, D., AND AGHA, G. CUTE: A concolic unit testing engine for C. In In 5th joint meeting of the European Software Engineering Conference and ACM Symposium on the Foundations of Software Engineering (ESEC/FSE 2005).
- [8] BOONSTOPPEL, P., CADAR, C., AND ENGLER, D. RWset: Attacking path explosion in constraint-based test generation. In Proceedings of Tools and Algorithms for the Construction and Analysis of Systems (TACAS 2008).
- [9] BRUMLEY, D., NEWSOME, J., SONG, D., WANG, H., AND JHA, S. Towards automatic generation of vulnerability-based signatures. In Proceedings of the 2006 IEEE Symposium on Security and Privacy (IEEE S&P 2006).
- [10] CADAR, C., AND ENGLER, D. Execution generated test cases: How to make systems code crash itself. In Proceedings of the 12th International SPIN Workshop on Model Checking of Software (SPIN 2005).
- [11] CADAR, C., GANESH, V., PAWLOWSKI, P., DILL, D., AND ENGLER, D. EXE: Automatically generating inputs of death. In Proceedings of the 13th ACM Conference on Computer and Communications Security (CCS 2006).
- [12] COSTA, M., CASTRO, M., ZHOU, L., ZHANG, L., AND PEINADO, M. Bouncer: Securing software by blocking bad input. In Proceedings of the 21th ACM Symposium on Operating Systems Principles (SOSP 2007).
- [13] COSTA, M., CROWCROFT, J., CASTRO, M., ROWSTRON, A., ZHOU, L., ZHANG, L., AND BARHAM, P. Vigilante: end-to-end containment of Internet worms. In Proceedings of the 20th ACM Symposium on Operating Systems Principles (SOSP 2005).
- [14] EMMI, M., MAJUMDAR, R., AND SEN, K. Dynamic test input generation for database applications. In International Symposium on Software Testing and Analysis (ISSTA 2007).
- [15] GODEFROID, P. Compositional dynamic test generation. In Proceedings of the 34th Symposium on Principles of Programming Languages (POPL 2007).
- [16] GODEFROID, P., KLARLUND, N., AND SEN, K. DART: Directed automated random testing. In Proceedings of the Conference on Programming Language Design and Implementation (PLDI 2005).
- [17] GODEFROID, P., LEVIN, M. Y., AND MOLNAR, D. Automated whitebox fuzz testing. In Proceedings of Network and Distributed Systems Security (NDSS 2008).
- [18] SCHWARZ, B., DEBRAY, S., AND ANDREWS, G. Disassembly of executable code revisited. In Working Conf. on Reverse Engineering, 2002.
- [19] DIILLIG, I., DILLIG, T., AND ALKEN, A. Sound, complete and scalable path-sensitive analysis. In Conf. on Programming Language Design and Implementation, 2008.
- [20] LAM, M. S., WHALEY, J., LIVSHITS, V. B., MARTIN, M. C., AVOTS, D., CARBIN, M., AND UNKEL, C. Context-sensitive program analysis as database queries. In Symp. on Principles of Database Systems, 2005.