
Selective Symbolic Execution

Anonymous, Technische Universität München

This paper describes the exemplary application of selective symbolic execution techniques for the analysis of a binary file in user mode. Goal of this study is to search the binary for possible privacy issues like unwanted leakage of personal data. Investigation will be done using S²E, a powerful platform for selective symbolic execution of large software systems.

1 Introduction

Frequently developers need to understand software systems. In a very simple case they just analyse their own code or test the interaction of own programs with other components or with the surrounding environment in general. Testing self-written programs conceptually permits the application of the whole arsenal of analysis techniques.

Things become interesting when analysis has to be performed without access to source code or documentation. Scenarios for this situation include the need to check proprietary third party software for interoperability on existing servers, performance, unwanted side effects, and much more. Security-critical environments additionally require reliable guarantees of the benignity of all employed software.

One mighty solution for such system analysis is the S²E platform developed at the Swiss Federal Institute of Technology in Lausanne (EPFL) [9]. Its goal is to provide a tool set for rapid development of analysis tools like performance profilers, bug finders, reverse engineering solutions and the like [12]. S²E combines several key characteristics:

1.) The ability to explore entire *families of execution paths* helps to obtain reliable information about the

target system. Abstracting from single-path exploration to sets of execution paths which share specific properties is vital for predictive analyses. This technique can for example prove the non-existence of critical corner cases which might be overlooked by other testing strategies.

2.) *In-vivo analysis*, meaning the analysis of a program within its real-world environment (libraries, kernel, drivers, etc.), facilitates extremely realistic and accurate results.

3.) Working directly on *binaries* further increases the degree of realism in system analyses, as it allows to include closed source modules into the investigation.

As an exemplary showcase for the power of the S²E platform and its underlying concepts this paper will perform a thorough analysis of a binary file in user mode. The scenario assumes that this program was found somewhere in the internet and claims to be a useful freeware tool. S²E shall help to investigate whether the binary compromises the user's privacy, for instance by leaking private data to the internet.

Chapter 2 explains the theoretical concepts of selective symbolic execution. Chapter 3 then introduces S²E, the platform which builds upon all techniques described before. Coming to the practical part, chapter 4 lays out the concrete analysis scenario and formulates research questions. Chapter 5 describes how S²E can be applied in this scenario, followed by an explanation of S²E results in chapter 6. Possible further research related to this topic is mentioned in chapter 7, together with some selected related work in chapter 8. Finally, chapter 9 summarises and concludes this paper.

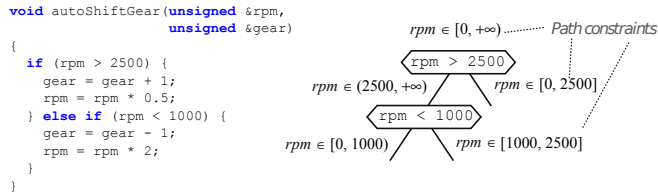


Figure 1: Execution tree with path constraints for the symbolic variable *rpm* [7]

2 Selective Symbolic Execution

Symbolic execution is an advanced analysis technique. Instead of concrete input (7, “string”, ...) symbolic execution uses symbolic values (λ , β , ...) when processing code. Assignments in the program path have impacts on these symbolic values. The integer calculation $x = x - 2$, for instance, would update the symbolic expression representing the input x to $\lambda - 2$. Conditional statements (*if <condition> then ... else ...*) fork program execution into two new paths. Both paths are then constrained by an additional condition, the ‘then’ branch with the if-condition and the ‘else’ branch with the negated if-condition respectively. Figure 1 conceptually shows how source code can be transformed to symbolic expressions.

Following this procedure results in a tree-like structure of constrained symbolic expressions. A constraint solver can now take all constraints along one execution path as input and find one concrete input (e.g., $\lambda = 5$) which would lead to the program following exactly this path. Such results greatly alleviate writing reproducible test cases [8].

On a technical level, symbolic execution engines save state information (program memory, constraint information, ...) in a custom data structure. Each conditional statement involving symbolic values results in a *fork* of the program state. The two newly created branches are completely independent and can therefore be processed in parallel.

But the exponential growth of conditionals soon reveals scaling problems of this forking strategy. Despite heavy research on optimisations mitigating this *path explosion* problem only relatively small programs (\cong thousands of lines of code) can be analysed symbolically [8].

Additionally, symbolic execution faces problems when the program under analysis *interacts with its environment*. If it calls a system library like *libc*, in theory the whole system stack including invoked libraries, operating system and drivers would have to be executed symbolically. Considering the path

explosion problem mentioned before, the resulting complexity makes such a profound analysis hardly feasible.

One way to solve this problem is to build abstract models of the program’s environment [6, 11]. However, due to the complexity of real-world systems, building a model of the entire system is both tedious and unnecessary - the user usually wants to analyse one single program and not the whole system [8].

In order to overcome typical problems of conventional symbolic execution, Chipounov et al. at EPFL developed the concept of **selective symbolic execution** (S²E) [8]. Based on a virtual execution platform S²E gives users the illusion of running the entire system symbolically. By limiting the scope of interest (i.e. which parts of the system should be executed symbolically), users can effectively restrain the path explosion problem. Program code within this defined scope is executed symbolically, whereas out-of-scope parts, which are irrelevant to the analysis, switch to concrete execution.

Definition of the scope of interest (what to execute symbolically) is highly flexible. Users may specify whole executables, code regions, or even single variables to be executed symbolically. Everything else will be treated concretely.

But since on a technical level symbolic and concrete execution are handled very differently - concrete code may run natively while symbolic instructions need to be emulated - switching back and forth these two modes is a major challenge. Hence one of the main contributions of the EPFL team around Chipounov is the transparent and consistent management of switching between symbolic and concrete execution modes.

Figure 2 depicts the **interplay of symbolic and concrete execution**. The illustration is based on a scenario where an application *App* is tested. A function *appFn* invokes the method *libFn* in a library *Lib*, which in turn calls a function *sysFn* in the kernel. Since we suspect a bug in *libFn*, we focus our analysis upon this function. Due to the path explosion problem, symbolically executing the entire system stack is not feasible. Hence only execution inside *libFn* follows this technique.

Concrete \rightarrow symbolic transition: When execution enters the function *libFn* it has to change from concrete into symbolic domain (grey areas). This is done by replacing concrete parameters in the method call with symbolic variables. The call *libFn*(10) becomes *libFn*(λ), optionally also with

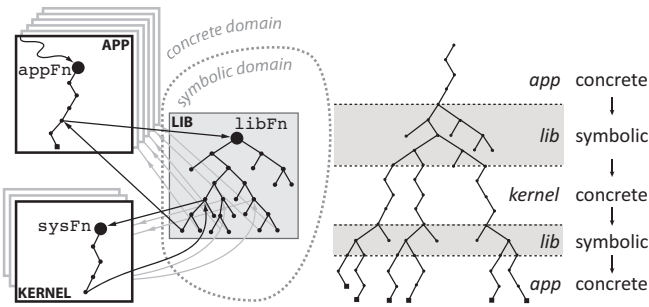


Figure 2: Selective symbolic execution: only paths inside a defined scope of interest (here: a library function `libFn`) are explored symbolically - the rest of the system stack runs concretely [7].

constraints: $libFn(\lambda \leq 15)$.

Besides the symbolic multi-path execution, S²E simultaneously also runs the function with its original concrete arguments. This is necessary in order to return a correct calculation result to `appFn` and thus keep the execution of `App` consistent.

Symbolic → concrete transition: Since the operating system is not focus of this analysis example, S²E has to switch from symbolic to concrete domain when `libFn` calls into the kernel. This is done by randomly picking a concrete value which fulfils all path constraints. If, for instance, the current path is constrained with $x \in]-\infty; 5]$, S²E might choose $x = 4$ and call `sysFn(4)`.

However, when `sysFn(4)` returns, `libFn` can no longer make any assumptions about any $x \neq 4$, because the behaviour of `sysFn` in those cases remains unclear. In order to preserve correctness, a new constraint $x = 4$ has to be added to the path¹. But imagine `libFn` being implemented as shown in figure 3; now the ‘then’ branch of the if-condition ($x < 0$) can never be reached. Chipounov calls this effect “overconstraining” [7] - it is a result of concretising x when leaving the symbolic domain. S²E tackles the problem by going back in the execution tree and forking an additional sub-tree. The new sub-tree now picks a different concrete value for x which allows to enter the previously unreachable ‘then’ branch.

The ability to flexibly switch between symbolic and concrete execution enables S²E to handle different **execution consistency models**. Chipounov defines six distinct consistency models. For a basic understanding, the following paragraph will describe two of them (see [7], p. 24ff for a detailed

¹Constraints added because of a symbolic → concrete transition are called ‘soft constraints’.

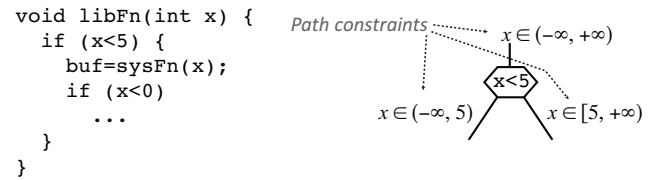


Figure 3: Excerpt from `libFn`’s execution tree [12]

explanation of all models).

The most simple model is called *strictly consistent concrete execution* (SC-CE). This is the model which is applied when testing or debugging programs either manually or with input from random test case generators (fuzzing). ‘Strict consistency’ describes the fact that every execution path can be mapped to concrete input values. This model never produces false positives, i.e. infeasible paths.

One entirely different strategy is *relaxed consistency overapproximate consistency* (RC-OC). In this model constraints defined in API contracts can be ignored. Under RC-OC, a call of the system function `write(fd, buf, count)` would also allow return values larger than `count`, even though that is infeasible (assuming that the C library works correctly). This behaviour can be particularly useful for reverse engineering or for finding vulnerabilities in a program. RC-OC is not consistent (pursues many infeasible paths) but complete (eventually finds every valid path through the system).

3 The S²E Platform

Based on the concepts described in the previous chapter, Chipounov and his team implemented the S²E platform, an open source framework for writing custom system analysis tools. S²E employs the theoretical concepts of selective symbolic execution by running the system under analysis in a virtual machine and treating code within the scope of interest as symbolic. These symbolic parts are translated into an intermediate representation (LLVM IR), while irrelevant instructions are directly passed to the host for native execution.

Technical backbone of S²E are the virtual machine hypervisor QEMU [3, 5], the symbolic execution engine KLEE [1, 6] and the LLVM compiler infrastructure [2, 10]. Figure 4 gives an overview of how these technologies are integrated into the S²E platform. The top of the picture depicts the software stack of the guest system (=the system under analysis), which is managed by QEMU. S²E is not restricted to user land applications, but also allows inspection

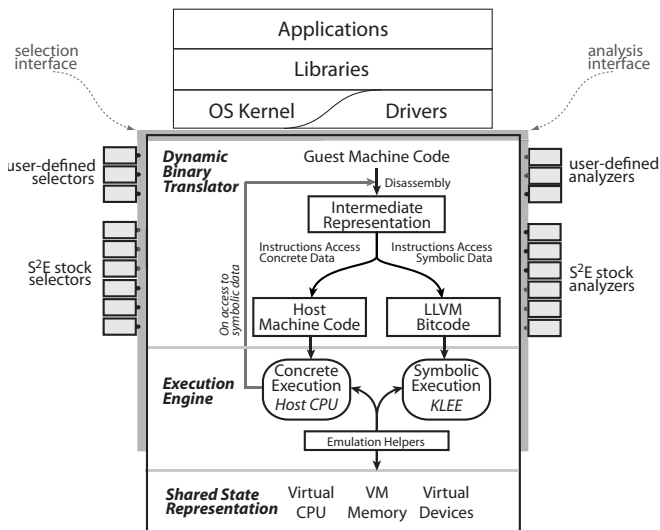


Figure 4: Architecture of the S²E platform [12]

on deeper levels (e.g., operating system functions).

For easier emulation, QEMU translates machine code of the guest system into an intermediate representation, called *microoperations*. S²E’s dynamic binary translator (DBT) splits the resulting microoperations into those that need to be explored symbolically and those which may run concretely. All concrete microoperations are directly converted into host instructions. Symbolic expressions, on the other hand, are prepared for being executed on the KLEE engine. This requires microoperations to be translated into the LLVM intermediate representation, called LLVM Bitcode in figure 4.

S²E’s execution engine, which is an extension to QEMU’s execution engine, now manages the operation of the platform. In an endless loop it asks the DBT for new guest code. Depending on the result, instructions can either be run straight on the host system or are fed into the KLEE symbolic execution engine.

In order to keep the mix of symbolic and concrete execution consistent, S²E stores state (VM CPU, memory, ...) centrally, by consolidating QEMU and KLEE data structures and managing them in a single machine state representation.

Users work with S²E by writing selection and analysis plugins or by simply configuring S²E’s standard plugins according to their needs. Plugins subscribe to system-wide events (e.g., *onInstrExecution*) and can perform logging/monitoring tasks or even manipulate the system state.

Configuration usually starts with defining what parts of the system to explore symbolically. This can for example be done with S²E’s selection plugin

CodeSelector, which restricts symbolic execution to a specified module or code region.

Standard analysis plugins allow users to find bugs (*WinBugCheck*), monitor memory (*MemoryChecker*), study performance characteristics (*PerformanceProfiler*) and much more (see [7], p. 50).

4 Project Idea and Research Questions

The practical part of this project strives to explore privacy issues in a sample binary.

In order to make life easier, many people use little freeware applications on a regular basis. But most of these programs are proprietary and have to be trusted without any knowledge of their functioning. Real malware (Trojan horses, spyware, ...) is usually detected rather quickly by anti-virus software and can often be blocked effectively. However, between unambiguous malware and thoroughly benign software many shades of grey can be found.

This work will focus on the scenario that an application (intentionally or unintentionally) leaks delicate private data without the user’s consent or knowledge.

Due to the difficulty of finding a real-world program which shows exactly this desired behaviour and also in general the complexity of real-world applications, the showcase described here bases on a little self-written program.

The scenario works as follows: The freeware tool *SuperTaxCalcPro* found in the internet claims to be handy for estimating your personal tax load. The tool announces to display a decent advertisement on every startup, which suggests a common and reasonably sound business model. It promises to treat all personal data confidentially, but since the application is closed source, of course this claim cannot be verified easily.

To make sure that *SuperTaxCalcPro* does not compromise the user’s privacy, it shall be analysed using selective symbolic execution techniques and the S²E platform.

For a general overview, often a good first step is to have a glance at the program’s assembly code. Since static code analysis of large programs is usually a very time-consuming task, that is not what we want to do here. Nevertheless, some useful information can be retrieved from the assembly code quite easily. A look at the list of invoked external

libraries shows that *SuperTaxCalcPro* communicates over network sockets, i.e. uses the C library functions `connect`, `write`, etc. That is of course not really surprising as the program relies on advertisements to the internet are always delicate and could potentially leak personal data. Therefore further analysis in this paper will focus on the following concrete questions:

1. When (under which conditions) are connections to the internet established?
2. What data is transferred in these connections?
3. Does a connection to the internet leak personal data?

These questions could of course be tackled via many different analysis techniques. Simply debugging the program should already lead to a rough understanding of the program's behaviour. Applying techniques like fuzzing could increase code coverage with the goal of not missing exotic or covert execution paths. This paper, however, will rely on analysis using selective symbolic execution as described in chapter 2.

Before proceeding with the actual analysis, a few notes on the general setup: Since S²E works best with x86 guest architectures, *SuperTaxCalcPro* was compiled on a x86 Debian system using standard GNU tools (g++). For the sake of convenience in this sample scenario the resulting application binary is neither overly optimised nor in any way obfuscated. As a consequence, the assembly code still contains meaningful variable and function names.

5 Implementation

Analysis using the S²E platform is usually done by writing a configuration script in the Lua programming language. As described in chapter 2, working with S²E can be divided into a code selection and the actual analysis part.

In the **code selection phase** S²E is configured to focus on the program *SuperTaxCalcPro* in user space only. The whole environment will be treated as a black box. Symbolic execution will only be applied inside the code of *SuperTaxCalcPro*.

The next step is to precisely specify which variables shall be treated symbolically. In general, as a start we want to make all user inputs symbolic. For

a closed source binary this can be done in several ways. A command line tool can be called from a wrapper program which hands over symbolic arguments instead of concrete ones. If user input is entered in a UI, a further look into the assembly code is necessary in order to find the memory locations where the corresponding variables are written. S²E can intercept execution of the specified memory addresses and replace them with symbolic values (`state.writeMemorySymb(...)`).

KLEE, the engine for symbolic execution, also requires some configuration. As path search strategy this analysis uses a depth first search.

```

1 08049415: <_Z9send_dataPc>:
2      .
3      .
4      8049451: e8 ba fa ff ff      call    8048f10 <socket@plt>
5      .
6      .
7      80494f2: e8 09 f9 ff ff      call    8048e00 <connect@plt>
8      80494f7: c1 e8 1f            shr     $0x1f,%eax
9      80494fa: 84 c0              test    %al,%al
10     80494fc: 74 0c              je      804950a
11     <_Z9send_dataPc+0xf5>
12     80494fe: c7 04 24 83 a0 04 08 movl    $0x804a083, (%esp)
13     8049505: e8 f9 fe ff ff      call    8049403 <_Z5errorPKc>
14     804950a: 8d 85 d8 fe ff ff      lea     -0x128(%ebp), %eax
15     8049510: 89 04 24            mov     %eax, (%esp)
16     8049513: e8 18 fa ff ff      call    8048f30 <strlen@plt>
17     8049518: 89 44 24 08         mov     %eax, 0x8(%esp)
18     804951c: 8d 85 d8 fe ff ff      lea     -0x128(%ebp), %eax
19     8049522: 89 44 24 04         mov     %eax, 0x4(%esp)
20     8049526: 8b 45 f0            mov     -0x10(%ebp), %eax
21     8049529: 89 04 24            mov     %eax, (%esp)
22     804952c: e8 4f f9 ff ff      call    8048e80 <write@plt>
23     .
24     .
25     8049564: 8d 85 d8 fe ff ff      lea     -0x128(%ebp), %eax
26     804956a: 89 44 24 04         mov     %eax, 0x4(%esp)
27     804956e: 8b 45 f0            mov     -0x10(%ebp), %eax
28     8049571: 89 04 24            mov     %eax, (%esp)
29     8049574: e8 67 f9 ff ff      call    8048ee0 <read@plt>
30     .

```

Listing 1: Relevant parts of the function `send_data` in assembly code. Interesting calls are highlighted.

Most logic in the **analysis part** relies on S²E's *Annotations* plugin. It allows fine monitoring and even fiddling with the execution state by annotating single instructions or functions in the program binary.

For the scenario in this paper all instructions which handle connections to the internet are of particular interest. The C library call `connect` helps to identify where the program is connecting to, and the corresponding `write` and `read` calls allow to find out what data is sent and received in this connection. Memory addresses of these relevant calls can be retrieved from the assembly code. Since all calls concerned with connecting to the internet seem to take place in a single function named `send_data`, calls to this function as a whole shall be monitored, too. Listing 1 shows relevant parts of the function `send_data` in assembly code.

```

1 pluginsConfig.Annotation = {
2   annotation_write = {
3     active=true,
4     module="prog",
5     address=0x804952c, -- write()
6     instructionAnnotation="do_write",
7     beforeInstruction=true,
8     switchInstructionToSymbolic=false
9   },
10  annotation_read = {
11    active=true,
12    module="prog",
13    address=0x8049574, -- read()
14    instructionAnnotation="do_read",
15    beforeInstruction=false,
16    switchInstructionToSymbolic=false
17  },
18  annotation_send_data = {
19    active=true,
20    module="prog",
21    address=0x8049415, -- send_data()
22    callAnnotation="call_ann",
23    paramcount = 1
24  }
25 }

```

Listing 2: Configuration of the Annotations plugin (part). Defines the instructions to be monitored and actions to trigger upon execution of these instructions. Note the link to the binary in listing 1 via memory addresses.

After providing S²E's Annotations plugin with all relevant memory addresses plus a little more configuration (shown in listing 2), we can now execute code every time the program runs into one of the interesting instructions. At this point we can ...

- 1.) log that instruction *X* was reached. Together with other S²E output this information shows which execution paths run into the *send_data* method and when they do so.

- 2.) save information about the executed instruction in the current plugin state. This allows to check whether there have been prior connections to the internet when the program runs into the next annotated instruction.

- 3.) read out parameters of the called function. The function *send_data*, for instance, is called with one parameter, which appears to be a memory location. Now S²E's analysis interfaces can be employed to dynamically find out what is written in the respective memory.

- 4.) read (or write) registers (*EAX*, *EBP*, ...). Thus S²E users can check and even manipulate the current execution state, similar to debug situations.

- 5.) switch variables (memory locations) from concrete to symbolic mode (or vice versa). This ex-

tremely precise control over the symbolic execution allows a fine-grained implementation of custom consistency models.

```

1 function call_ann (state, curPlgState)
2   if curPlgState.isCall() then
3     arg0 = state.readParameter(0);
4     print ("Calling function with with
5       arg0=" .. ("%x"):format(arg0));
6     -- print_mem dereferences the pointer
7       arg0 and converts the memory
8       behind that address into a string.
9     print_mem(state, curPlgState, arg0);
10    -- saving info about the no. of
11      connections in this state so far
12    local no = curPlgState.getValue("no");
13    if no == nil then
14      no = 0
15    end
16    no = no + 1;
17    curPlgState.setValue("no", no);
18    print ("\n\tNo. of connections on
19      this path so far: " .. no);
20  else
21    -- Called when the function send_data
22      returns. The return value can be
23      read from the EAX register.
24    local eax = state.readRegister("eax");
25    print ("EAX: " .. ("%x"):format(eax));
26  end
27 end

```

Listing 3: Lua function executed upon every call of the function *send_data*(). See lines 18 - 24 in listing 2 for the registration of *call_ann*.

Listing 3, for instance, shows the Lua function which is triggered on every execution of the assembly function *send_data* (see listing 2 for the registration of this function and its coupling to a specific memory address). Because *send_data* was registered as a function call in listing 2, S²E's interfaces facilitate the extraction of all parameters via *state.readParameter()*. A sample test run shows that *arg0* is clearly a memory address, so the next step must be to examine what data is stored at this location. The method *print_mem* does just that and additionally brings the contents into a human-readable format. Then the Lua script stores information that this function was called in the current execution state. Doing so, the next call will notice that the method has been executed before, increment the counter and write this fact into a log file.

The *else* branch will be called when the function returns. Following common calling conventions, we expect return values to be transmitted via the register *EAX* (and *EDX* if bigger than four bytes), hence *EAX* could also hold interesting information.

Unfortunately, this does not seem to be the case as *EAX* is always 0 except in some error cases. So *send_data* is probably just returning a status integer value.

In addition to the *Annotations* plugin, analysis in this paper also employs the *TestCaseGenerator*. For each execution path, it finds concrete inputs for all symbolic variables. Those will be printed upon termination of the path and, apart from helping to understand the program, may serve as input for further testing.

Naturally, an important step for understanding symbolic execution of *SuperTaxCalcPro* is to log and later interpret output of KLEE, the symbolic execution engine used in S²E. The behaviour of KLEE can also be controlled via S²E's Lua configuration file.

6 Analysis of S²E Output

Each execution of S²E creates a folder for analysis output. Apart from the general log files, where all plugins write important messages, some plugins also create separate files. For instance one file contains performance characteristics of the test run, another one LLVM bytecode of all analysed programs. The exact numbers and types of output files depend on which plugins are defined in S²E's configuration Lua script.

The most important information about S²E's general execution is accumulated in *messages.txt*. Backbone of this file is information about all execution paths, at which memory addresses they were forked, how they are constrained, when and why they were terminated, and much more. It is also the place where we saved all custom log messages as defined in the previous chapter to.

Thanks to the clear structure of *messages.txt* it can be used as input for a custom Python script. This script was written in order to visualise the **tree of execution paths** graphically, which serves as a great help with understanding path forking behaviour. Figure 5 shows the vital part of the tree of execution paths of *SuperTaxCalcPro*.

Each box represents an execution state, lines pointing to another state symbolise a fork of the execution state. The hexadecimal number at the origin of each edge is the memory address in which the state was forked. For the sake of clarity I also manually added the function in which each memory address occurs.

Edges are labelled with constraints that need to be respected in the corresponding execution state. KLEE handles constraints in Extended Backus-Naur Form syntax (e.g., "(Not (Slt (w32 ox20a1) (ReadLSB w32 oxo vo_income_o)))" for " $income \geq 8353$ "). For simplicity all constraints were also manipulated manually and transformed into a more readable format.

The text below each leaf shows the output of the *TestCaseGenerator* plugin for this execution path. Upon termination of a state it finds concrete input examples for the two symbolic variables *income* and *taxcat*.

Before coming to the actual research objectives of this paper, namely finding privacy problems, a quick interpretation of the bare tree shown in figure 5 helps to gain a common understanding of *SuperTaxCalcPro*.

Program execution largely relies on two integer variables, *income* and the tax category *taxcat*. Most forks in the assembly function *main()* seem to be basic checks if these two variables contain valid data (e.g., [*taxcat* > 6]). This assumption arises from the fact that many of these forks are leaves which exist only a very short time and consume hardly any memory (can be checked via the *MemoryChecker* plugin). These paths are probably error cases that terminate the program.

Important logic seems to take place in a function called *calc_tax()*. Many forks originate in this method and both symbolic variables play a role in it. Important in this case means 'important with respect to being relevant to the symbolic execution of the two focal variables *income* and *taxcat*'. Paths with the constraints [*income* < 13468] and [*income* < 52880] seem to be less important since they are leaves in the tree. The two paths [$52880 \leq income < 250729$] and [$income \geq 25729$], however, entail further checks of *taxcat*, the other symbolic variable. Summarising the general findings so far, *SuperTaxCalcPro* checks the user's income and tax category and takes special actions if the income is higher than 52880 euros.

Unfortunately, due to the problem with incomplete test results (described at the end of this chapter), the tree in figure 5 does not show all execution states. While many states (900+) were manually left out due to insignificance, many more states are missing in S²E's output despite potentially being relevant. For example, we cannot say what happens if *income* is below or equal 8353 (the complementary event of the switch from state 3 to state 4).

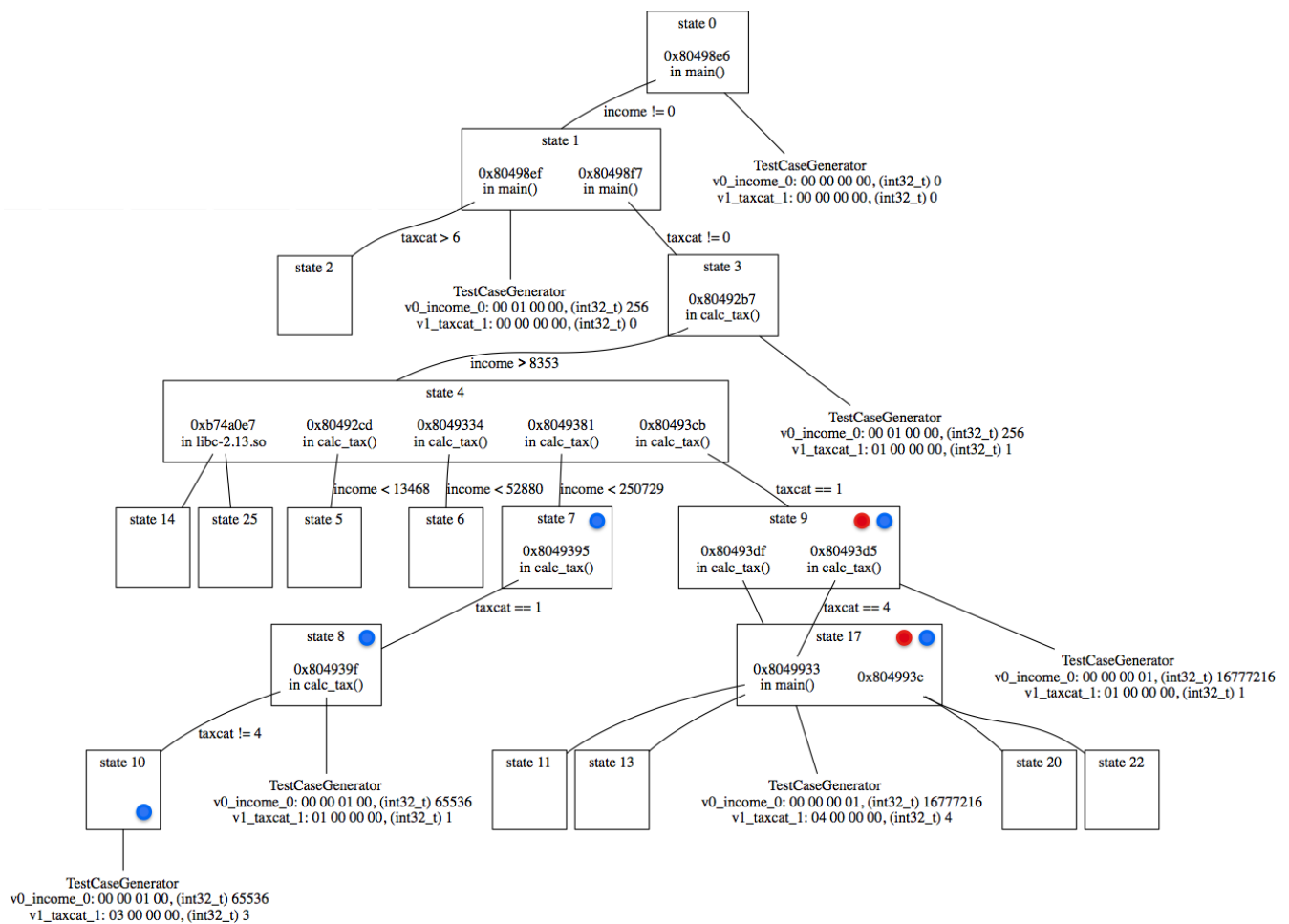


Figure 5: Tree visualisation for understanding path forking behaviour in the analysis of SuperTaxCalcPro. Shows at which memory addresses S²E forked new execution states. Path constraints are printed at each state transition. For most states, the TestCaseGenerator found meaningful concrete example values for the two symbolic variables `income` and `taxcat`. Blue dots symbolise that in this state a “donation” message was transferred to a server. Red dots mark states which leak the complete set of user inputs.

After this general inspection of *SuperTaxCalcPro*'s behaviour, we can finally shift the focus towards the research objectives regarding privacy issues. As described in the previous chapter, S²E was programmed to monitor all instructions related to connections to the internet. Output of these inspection activities have to be extracted from the S²E's main log file. Based on the output resulting from annotating the function *send_data*, which includes reading out the parameter the method is called with (see listing 3), messages sent to the internet can be classified into three categories:

1. A command "load_ad", which leads to a server response containing an advertisement for the user.
2. A command "donate", which entails a static message from the server prompting the user to donate to the poor programmer. Output messages 1.) and 2.) do not contain any further data apart from the command.
3. A message depending on both symbolic variables *income* and *taxcat*. It always starts with the keyword "found_hit", followed by all data the user has entered. Message fields are separated by a colon and stick to the order 'user name', 'birth date', 'income' and 'tax category'. Example: "found_hit:name:1.1.77:16777216:1". Response from the server is a static status message.

While message type 1.) could be expected (*SuperTaxCalcPro* informed us that it shows advertisements), calls for donation (type 2.)) were not announced before. Still, since "donate" messages here do not contain any user-specific data, at this point of the analysis they can be considered innocuous with regard to privacy issues.

Internet connections of type 3.), however, undoubtedly impair the user's privacy. They send the whole set of confidential user data to a server in the internet. Hence a clear answer to research question three is: yes, *SuperTaxCalcPro* does leak personal data of the user!

Note that even if *SuperTaxCalcPro* would try to disguise its actions by encrypting the data it sends to the server, S²E would still find out which symbolic variables the encrypted message depends on. Thus it should be possible to find out what data is sent even without thorough analysis of the encryption mechanisms.

Now that all different connections to the internet are unveiled, we want to take advantage of S²E's

ability to pursue all possible execution paths and find out in which cases *SuperTaxCalcPro* connects to the internet.

For this purpose the execution tree in figure 5 was enriched with the information in which execution state which connections to the internet took place.

Blue dots in the graph mark that this state triggered a donation call (message of type 2.)). A look at path constraints in these states tells that the common criteria here must be an *income* either between 52880 and 250729 or even above 250729 euros. Apparently the programmer assumed wealthy people to be more generous than those with a low income. Despite the fact that these type 2.) calls do not transfer any user data to the server, they may potentially be a delicate matter since merely triggering them implicitly reveals information about the current user's income to the server.

Red dots in the tree visualise a message leaking data in this state (message type 3.)). Since all affected states originate at the fork in memory address 0x80493cb, one condition for type 3.) connections must be an *income* above 250729 euros. Results of the *TestCaseGenerator* plugin for the states 9 and 17 confirm this observation. Additionally, path constraints require *taxcat* to be either 1 or 4 in order to trigger this type of connection. Obviously the evil programmer of *SuperTaxCalcPro* is only interested in data of users with a very high income and tax category one or four.

The advertisement was loaded in every state, so this connection was not marked in the execution tree. *SuperTaxCalcPro*'s source code apparently triggers the call for an advertisement before checking anything related to the symbolic variables *income* and *taxcat*.

Since S²E eventually finds all possible execution paths of a program (if configured to do so), we can be sure that we did not miss any message type or occasion where connections are established. When simply debugging *SuperTaxCalcPro* it is quite likely that we do not hit the two cases where the application leaks data and thus remain ignorant of the severe privacy problems.

Summing up, we can now answer all research questions formulated in chapter 4:

1. When (under which conditions) are connections to the internet established?
 ⇒ *SuperTaxCalcPro* always asks its server for an advertisement.
 ⇒ Donation messages are sent whenever [*income* ≥ 52880].

⇒ User data is transferred if $[income \geq 250729 \ \&\& \ (taxcat == 1 \ || \ taxcat == 4)]$.

2. What data is transferred in these connections?
⇒ Type 1.) and type 2.) connections only send a short command string, while type 3.) messages contain all personal data the user has entered.
3. Does a connection to the internet leak personal data?
⇒ Yes, in type 3.) communications *SuperTaxCalcPro* reveals the user's name, his birth date, income and tax category to its server.

Unfortunately, concrete test runs in S²E were also accompanied with several **problems**:

1.) Even several hours of execution did not suffice to complete the analysis process. Every execution had to be terminated prematurely. Hence analysis output (list of states, etc.) was never complete. Probable causes of this problem are the following:

- Most importantly, the general path explosion as explained in the theoretical part of this paper.
- Following many execution paths that are irrelevant to the analysis objectives.
- Execution of S²E on a ordinary laptop and not on a powerful server.
- Slow communication of the analysed application with a server which simultaneously ran inside the same virtual machine.
- Inefficient configuration of S²E: Potentially unnecessary states being forked, too much debug and log output, ...

The longest analysis run was terminated manually after about six hours. It pursued over 1000 execution states. This run produced 164 MB of analysis output, the important pure text file logs being 18 MB big.

2.) Due to the extreme duration of analysis runs the original idea to treat all user input as symbolic had to be dismissed. Instead, only two variables (*income* and *taxcat*) are made symbolic now. This may be seen as minor cheating since I knew which variables would turn out to be relevant. A realistic setting will potentially require to run the same analysis with all possible combinations of symbolic variables in order to find a good set of symbolic variables which lead to meaningful analysis results.

3.) The initial idea to pursue states in a depth-first search manner turned out to be unwise. The parts of the execution tree which are vital for a grasp of the program under analysis tend to be close to the root of the tree. Keeping in mind that all S²E runs had to be terminated prematurely, conducting a depth-first search reveals a multitude of extremely specialised branches and corner cases while neglecting fundamental paths of the tree.

4.) The problem of incomplete test results also entailed complications while constructing the execution tree, which of course never was consistent. For this reason and due to the fact that most states are irrelevant for the goals in this paper anyway, the tree in figure 5 shows a radically trimmed version of the original execution tree.

However, despite all these problems, the quality of analysis output in the scenario of this paper was sufficient to produce respectable and sound results.

7 Outlook

This chapter delineates a few suggestions for further interesting research topics related to the project of this paper, starting with concrete ideas for further analysis of the scenario defined in chapter 4.

Due to the path explosion problems described above it would obviously be interesting to run S²E on a more powerful machine. This would allow a more comprehensive observation strategy, like using more symbolic variables.

Very fine-grained scripting of S²E's configuration file or even interfering in S²E's C++ source code could handle state forking and symbolic-concrete switching more intelligently. Out of the box, S²E of course does not know our analysis goals and pursues every state. Adapting this internal behaviour to personal needs should reduce execution time and effort significantly.

During the analysis done for this paper, the server responding to *SuperTaxCalcPro*'s messages was available as binary and always running on the guest VM. This is certainly a very unlikely situation when analysing real-world programs. The server might not be available or respond with much latency. It could potentially block the overload of requests triggered by our test execution, or we might not want to reveal our analysis efforts.

In all of these cases symbolic execution must become independent of the server's presence. Besides,

SuperTaxCalcPro's execution tree currently does not depend on server responses. If the application contained switches like *'if srv_response == x then do ...'*, then the consistency model applied in chapters 5 and 6 would fail.

Such a scenario would require S²E execution to be configured for using a *overapproximate consistency* (RC-OC) or *control flow graph consistency* (RC-CC) model (see [7], p. 26ff). In the *SuperTaxCalcPro* example, this would mean that responses from the server have to be intercepted and replaced with unconstrained symbolic variables. Practically speaking, the whole function *send_data* would have to be replaced with a symbolic variable. As a consequence, S²E would of course produce lots of states based on completely absurd potential server responses. We simply do not know how the server responds to a certain client message nor how a valid server response looks like. But this is the only way to make sure that S²E eventually also finds all valid server responses in the absence of any knowledge of how the server works.

Moving away from the the concrete example program used this paper, a really interesting research topic could be to apply similar analysis techniques to real-world applications. Maybe an analysis procedure similar to the one described in chapters 5 and 6 might work for the discovery of intended or unwanted (~bugs) cases of data leakage in publicly available software.

On a general level S²E may also be helpful for analyses of malware samples. This would of course be a challenging task and require different project objectives and a much more advanced test setup, but could nonetheless be feasible considering S²E's holistic approach.

8 Related Work

Banabic et al. do bla... [4]

BitBlaze (see rel work in chip11s2e)

9 Conclusion

This paper described an approach for finding data leaks in a binary file using selective symbolic execution techniques provided by the S²E platform. The investigated binary file, a self-written application called *SuperTaxCalcPro*, claimed to be a useful, ad-supported freeware tool, but under certain circumstances turned out to transmit all confidential

user data to a server in the internet.

S²E proved to be a mighty platform for universal testing of binary programs and even entire systems. With a little practice, its code selection and analysis interfaces can be controlled conveniently via Lua configuration scripts. After adequate configuration, S²E's standard analysis output already gives solid insights into the functioning of the investigated binary. Analysis output can be enhanced by writing custom Lua code.

But even under the laboratory settings employed here, with a very small and simple binary, path explosion of the symbolic execution turned out to be a big problem concerning time and the resources it takes to explore all execution states. Investigation of real-world applications would require a more intelligent handling of S²E's path forking behaviour, with the goal of pursuing only states which are relevant for the analysis.

References

- [1] KLEE, <https://klee.github.io>.
- [2] LLVM, <http://llvm.org>.
- [3] QEMU, <http://qemu.org>.
- [4] Radu Banabic, George Candea, and Rachid Guerraoui, *Finding Trojan Message Vulnerabilities in Distributed Systems*, Proceedings of the 19th international conference on Architectural support for programming languages and operating systems, ACM, 2014, pp. 113–126.
- [5] Fabrice Bellard, *QEMU, a Fast and Portable Dynamic Translator*, USENIX Annual Technical Conference, FREENIX Track, 2005, pp. 41–46.
- [6] Cristian Cadar, Daniel Dunbar, and Dawson R Engler, *KLEE: Unassisted and Automatic Generation of High-Coverage Tests for Complex Systems Programs*, OSDI, vol. 8, 2008, pp. 209–224.
- [7] Vitaly Chipounov, *S2E: A Platform for In-Vivo Multi-Path Analysis of Software Systems*, Ph.D. thesis, École Polytechnique Fédérale de Lausanne (EPFL), 2014.
- [8] Vitaly Chipounov, Vlad Georgescu, Cristian Zamfir, and George Candea, *Selective Symbolic Execution*, 5th Workshop on Hot Topics in System Dependability (HotDep), 2009.

- [9] Vitaly Chipounov, Volodymyr Kuznetsov, and George Candea, *S2E: A Platform for In-Vivo Multi-Path Analysis of Software Systems*, 16th Intl. Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS), 2011.
- [10] Chris Lattner and Vikram Adve, *LLVM: A Compilation Framework for Lifelong Program Analysis & Transformation*, Code Generation and Optimization, 2004. CGO 2004. International Symposium on, IEEE, 2004, pp. 75–86.
- [11] Corina Pasareanu, Peter C Mehlitz, David Bushnell, Karen Gundy-Burlet, Michael Lowry, Suzette Person, and Mark Pape, *Combining Unit-Level Symbolic Execution and System-Level Concrete Execution for Testing NASA Software*, Proceedings of the 2008 international symposium on Software testing and analysis, ACM, 2008, pp. 15–26.
- [12] Vitaly Chipounov, Volodymyr Kuznetsov, and George Candea, *The S2E Platform: Design, Implementation, and Applications*, ACM Transactions on Computer Systems (TOCS) **30** (2012).