

# AUTOMATICALLY LOCATING ARM INSTRUCTIONS DEVIATION BETWEEN REAL DEVICES AND CPU EMULATORS

A PREPRINT

Muhui Jiang<sup>1,2</sup>, Tianyi Xu<sup>1</sup>, Yajin Zhou<sup>\*1</sup>, Yufeng Hu<sup>1</sup>, Ming Zhong<sup>1</sup>, Lei Wu<sup>1</sup>, Xiapu Luo<sup>2</sup>, and Kui Ren<sup>1</sup>

<sup>1</sup>*Zhejiang University*

<sup>2</sup>*The Hong Kong Polytechnic University*

August 13, 2021

## ABSTRACT

Emulator is widely used to build dynamic analysis frameworks due to its fine-grained tracing capability, full system monitoring functionality, and scalability of running on different operating systems and architectures. However, whether the emulator is consistent with real devices is unknown. To understand this problem, we aim to automatically locate inconsistent instructions, which behave differently between emulators and real devices.

We target ARM architecture, which provides machine readable specification. Based on the specification, we propose a test case generator by designing and implementing the first symbolic execution engine for ARM architecture specification language (ASL). We generate 2,774,649 representative instruction streams and conduct differential testing with these instruction streams between four ARM real devices in different architecture versions (i.e., ARMv5, ARMv6, ARMv7-a, and ARMv8-a) and the state-of-the-art emulators (i.e., QEMU). We locate 155,642 inconsistent instruction streams, which cover 30% of all instruction encodings and 47.8% of the instructions. We find undefined implementation in ARM manual and implementation bugs of QEMU are the major causes of inconsistencies. Furthermore, we discover four QEMU bugs, which are confirmed and patched by the developers, covering 13 instruction encodings including the most commonly used ones (e.g., STR, BLX). With the inconsistent instructions, we build three security applications and demonstrate the capability of these instructions on detecting emulators, anti-emulation, and anti-fuzzing.

## 1 Introduction

CPU emulator is a powerful tool as it provides fundamental functionalities (e.g., tracing, record and replay) for the dynamic analysis. Though hardware-based tracing techniques exist, they have limitations compared with software emulation. For example, ARM ETM has limited Embedded Trace Buffer (ETB). The size of ETB of the Juno Development Board is 64KB<sup>2</sup> [1]. On the contrary, software emulation is capable of tracing the whole program, provides user-friendly APIs for runtime instrumentation, and can run on multiple operating systems (e.g., Windows and Linux) and host machines in different architectures. Nevertheless, software emulation complements the hardware-based tracing and provides rich functionalities that dynamic analysis systems can build upon.

Indeed, many dynamic analysis frameworks [14, 18, 44, 9, 40, 10, 12, 21, 11, 16, 31, 26, 22] are built based on the state-of-the-art CPU emulator, i.e., QEMU, to conduct malware analysis, live-patching, crash analysis and etc. Meanwhile, many fuzzing tools utilize CPU emulators to fuzz binaries, e.g., the QEMU mode of AFL [2], FirmAFL [45], P2IM [15], HALucinator [13] and TriforceAFL [7].

The wide adoption of software emulation usually has an implicit assumption that the execution result of an instruction on the CPU emulator and the real device is identical, thus running a program on the CPU emulator can reflect the result

<sup>\*</sup>Corresponding author (yajin.zhou@zju.edu.cn).

<sup>2</sup>The ETB size of different SoCs may be different. However, it's usually limited due to the chip cost and size.

on the real hardware. *However, whether this assumption really holds in reality is unknown.* In fact, the execution result could be different (as shown in our work), either because the CPU emulator has bugs or because it uses a different implementation from the real device. These differences impede the reliability of emulator-based dynamic analysis. For instance, the malware can abuse the differences to protect the malicious behaviors from being analyzed in the emulator [37, 20, 19, 30].

In this work, we aim to automatically locate inconsistent instructions between real devices and the CPU emulator for the ARM architecture. If an instruction behaves differently between them, then it is an inconsistent instruction. Although previous research [32, 33, 35, 34] provides valuable insights, they are limited to the x86/x64 architecture and cannot be directly applied to the ARM architecture. Our work leverages the differential testing [36] for the purpose. Specifically, we provide the same instruction stream<sup>3</sup> to both the real device and a CPU emulator, and compare the execution result to check whether it is an inconsistent one.

Though the basic idea is straightforward, it faces the following two challenges. *First*, the ARM architecture has multiple versions (e.g., ARM v5, v6, v7 and v8), different register widths (16 bits or 32 bits) and instruction sets. Besides, it has mixed instruction modes (ARM, Thumb-1 and Thumb-2). Thus, how to generate effective test cases, i.e. instruction streams that cover previously mentioned architecture variants, while at the same time generating only necessary test cases to save the time cost, is the first challenge. Notice that if we naively enumerate 32-bit instruction streams, the number of test cases would be  $2^{32}$ , which is inefficient, if not possible, to be evaluated. Meanwhile, randomly generated instruction streams are not representative and many instructions are not covered (Section 4.1). *Second*, for each test case, we should provide a deterministic environment to execute the single instruction stream and automatically compare the result after the execution. This requires us to set up the same context (with CPU registers and memory regions) before the execution and compare the context afterwards.

Our system solves the challenges with the following two key techniques.

**Syntax and semantics aware test case generation** To generate representative instruction streams, we propose a syntax and semantics aware test case generation methodology. Each ARM instruction consists of several *encoding schemas*, which is called *instruction encodings* in this paper, that define the instruction’s structure (syntax). Each encoding schema maps to one decoding and execution logic that defines the instruction semantics. The encoding schema shows which parts of an instruction are constants and which parts can be mutated (Figure 2(a)). The non-constant parts of an instruction are called *encoding symbols* in this paper. The decoding and execution logic is expressed in the ARM’s Architecture Specific Language (ASL) [38]. We call it the *ASL code* in this paper (Figure 2(b) and (c)). The ASL code executes based on the concrete values of the encoding symbols. For instance, if the concrete value of the encoding symbol *W* (the eighth bit of *STR (immediate)* instruction) is 1, then the new address will be written back into the destination register *Rn* (line 4 of Figure 2(c)).

Specifically, during the test case generation, we first take the syntax-aware strategy. For each encoding symbol, we mutate it based on pre-defined rules. For instance, for the immediate value symbol, the values in the mutation set cover the maximum value, the minimum value and a fixed number of random values. This strategy generates syntactically correct instructions.

We further take a semantics-aware strategy to generate more instruction streams. That’s because the previous strategy may only cover limited instruction semantics as different encoding symbol values can result in different decoding and executing behaviors (Section 3.1). To this end, we extract the constraints in ASL code of decoding and executing. We solve the constraints and their negations by designing and implementing the first symbolic execution engine for ASL to find the satisfied values of the encoding symbols. By doing so, the generated test cases can cover different semantics of an instruction.

**Deterministic differential testing engine** Our differential testing engine uses the generated test cases as inputs. To get a deterministic testing result, we provide the same context when executing an instruction stream on a real CPU and an emulator. Besides, an instruction stream cannot be directly loaded and executed by the emulator, we carefully design a template binary that converts one instruction stream to a testing binary by inserting the prologue and epilogue instructions. The prologue instructions aim to set the execution environment while the epilogue instructions will dump the execution result for comparison to check whether the testing instruction stream is an inconsistent one.

We have implemented a prototype system called INSDet. Our test case generator generated 2,774,649 instruction streams that cover all the 1,998 ARM instruction encodings from 1,070 instructions in four instruction sets (i.e., A64, A32, T32, and T16).

---

<sup>3</sup>In this paper, instruction and instruction stream represent different meanings. For example, we call *STR (immediate)* an instruction. We call the concrete bytecode (i.e., 0xf84f0ddd) an instruction stream. See Section 2.1

On the contrary, the same number of randomly generated instruction streams can only cover 51.4% Instructions. This result shows the sufficiency of our test case generator.

We then feed these test cases into our differential testing engine. By comparing the result between the state-of-the-art emulator (i.e., QEMU) and real devices with four architecture versions (ARMv5, ARMv6, ARMv7-a, and ARMv8-a), our system detected 155,642 inconsistent instruction streams. Furthermore, these inconsistent instruction streams cover 47.8% of the instructions.

We then explore the root causes of them. It turns out that implementation bugs of QEMU and the undefined implementation in the ARM manual (i.e., the instruction does not have a well-defined behavior) are the major causes. We discovered four implementation bugs of QEMU and all of them have been confirmed by developers. These bugs influence 13 instruction encodings, including commonly used instructions, e.g., BLX, STR.

To show the usage of our findings, we further build three applications, i.e., emulator detection, anti-emulation and anti-fuzzing. By (ab)using inconsistent instructions, a program can successfully detect the existence of the CPU emulator and prevent the malicious behavior from being monitored by the dynamic analysis framework based on QEMU. Besides, the coverage of the program being fuzzed inside an emulator can be highly decreased. Note that, we only use these applications to demonstrate the usage scenarios of our findings. There may exist other applications, and we do not claim the contribution of them in this paper.

Our work makes the following main contributions.

**New test case generator** We propose a test case generator by introducing the first symbolic execution engine for ARM ASL code. It can generate representative instruction streams that sufficiently cover different instructions (encodings) and semantics.

**New prototype system** We implement a prototype system named INSDet that consists of a test case generator and a differential testing engine. Our experiments showed INSDet can automatically locate inconsistent instructions.

**New findings** We explore and report the root cause of the inconsistent instructions. Implementation bugs of QEMU and undefined implementation in ARM manual are the major causes. Furthermore, four bugs have been discovered and confirmed by QEMU developers. Some of them influence commonly used instructions (e.g., STR, BLX).

We will release generated test cases and the source code of our system to engage the community.

## 2 Background

### 2.1 Terms

For better illustration and avoid the potential confusion. We give detailed definition towards the following terms used in this paper.

**Instruction** Instruction denotes the category of ARM instructions in terms of functionality, which is usually represented by its name in ARM manual. For example, STR (immediate) is an instruction, which aims to store a word from a register to memory.

**Instruction Encoding** Instruction encoding refers to the encoding schemas for each instruction. We also call it encoding diagram in this paper. One instruction can have several encoding schemas.

**Instruction Stream** Instruction stream refers to the bytecode of an instruction. For example, 0xf84f0ddd, which meets one of the encoding schema of instruction STR (immediate). We call 0xf84f0ddd an instruction stream.

### 2.2 ARM Instruction and Instruction Encoding

Processor specification is important as it can verify the implementation of hardware, compilers, emulators, etc. To formalize the specification, ARM introduced the Architecture specification language (ASL) [38], which is machine-readable and executable.

ARM instructions usually have a fixed length (16 bits or 32 bits). According to ARM manual, one instruction may consist of several different instruction encodings, which describe the instruction structure (syntax). Our system generates the instruction streams that cover all the instruction encodings (which cover all instructions.) Specifically, the instruction encoding describes which parts of the instruction are constant and which parts are not. Each instruction encoding is further described with specific decoding and executing logic. The decoding and executing logic (expressed in ASL) defines the semantics of the instruction.

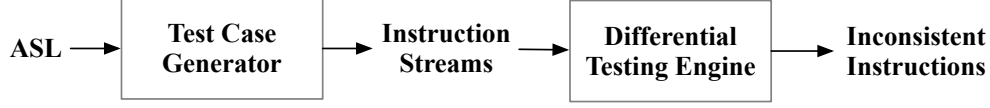
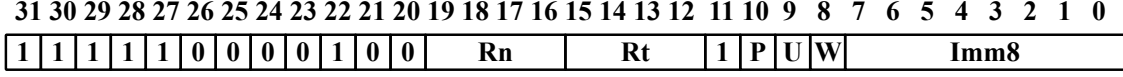


Figure 1: The work flow of our system



(a) The encoding schema of the STR (immediate) instruction in Thumb-2 mode.

```

1  if Rn == '1111' || (P == '0' && W == '0') then UNDEFINED;
2  t = UInt(Rt);
3  n = UInt(Rn);
4  imm32 = ZeroExtend(imm8, 32);
5  index = (P == '1');
6  add = (U == '1');
7  wback = (W == '1');
8  if t == 15 || (wback && n == t) then UNPREDICTABLE;
  
```

(b) The ASL code for decoding the instruction.

```

1  offset_addr = if add then (R[n] + imm32) else (R[n] - imm32);
2  address = if index then offset_addr else R[n];
3  MemU[address, 4] = R[t];
4  if wback then R[n] = offset_addr;
  
```

(c) The ASL code for executing the instruction.

Figure 2: A motivating example.

## 2.3 Instruction Decoding in QEMU

QEMU is the state-of-the-art CPU emulator that supports multiple CPU architectures. When executing an instruction stream, it needs to decode the instruction stream. QEMU adopts a two-stage decoding schema. In the first stage, it matches an instruction stream with pre-defined patterns, each of them represent multiple instructions. Then it distinguishes each instruction encoding based on the concrete value of the instruction. For instance, QEMU groups VLD4, VLD3, VLD2, and VLD1 instruction into one group (with one common pattern) and then identifies them inside the instruction decoding routine. If no instruction pattern can be found or further decoding routine cannot recognize an instruction stream, the SIGILL signal will be raised for the user mode emulation of QEMU.

## 3 Design and Implementation

Figure 1 shows the workflow of INSDet, which consists of a test case generator and a differential testing engine. First, the test case generator retrieves the ASL code to generate the test cases (Section 3.2). Then, the differential testing engine receives the generated test cases and conducts differential testing between the state-of-the-art emulator (i.e., QEMU) and real devices (Section 3.3). The instructions that can result in different behaviors are reported as inconsistent instructions. We further analyze the identified inconsistent instructions to understand the root cause of them and how they can be (ab)used.

In the following, we first use an inconsistent instruction detected by our system as a motivating example (Section 3.1), and then elaborate the test case generator and the differential testing engine in Section 3.2 and Section 3.3, respectively.

### 3.1 A Motivating Example

#### 3.1.1 The encoding schema and semantics of the STR (immediate) instruction

Figure 2 shows one of the encoding schema of instruction STR (immediate) and the corresponding ASL code for decoding and execution logic. According to the encoding schema in Figure 2a, the value is constant (i.e., 111110000100) for offset [31:20]. The encoding symbol Rn and Rt represent the addressing register and the source register, respectively. The last eight bits ([7:0]) represents a symbol value named imm8 that will be used as the offset.

Figure 2b shows the ASL code of the decoding logic for the encoding schema. Note that the ASL code is simplified for presentation. The complete code can be found in ARM official site [3].

- The ASL code at Line 1 checks the value of  $R_n$ ,  $P$ , and  $W$ . If the conditions are satisfied (or constraints are met), the instruction stream will be treated as an UNDEFINED one. Consequently, a SIGILL signal will be raised in QEMU user mode emulation when an UNDEFINED instruction stream is executed.
- In line 2 and 3, the symbol  $R_t$  and  $R_n$  will be converted to unsigned integer  $t$  and  $n$ , respectively. Similarly, the symbol  $imm8$  will be extended into a 32-bit integer  $imm32$ . In line 5, 6, and 7, symbol  $index$ ,  $add$ , and  $wback$  will be assigned according to the value of  $P$ ,  $U$ , and  $W$ , respectively.
- In line 8, the symbol  $t$ ,  $wback$ , and  $n$  will be checked. If the constraint of each condition is met, the instruction stream should be treated as an UNPREDICTABLE one. According to ARM’s manual, the behavior of an UNPREDICTABLE instruction stream is not defined. The CPU processor vendors and the emulator developers can choose an implementation that they think it’s proper.

Similarly, Figure 2c shows the ASL code for the execution logic of the instruction. The ASL code in Figure 2b and Figure 2c defines the semantics of the instruction.

### 3.1.2 Test case generation

By analyzing the encoding schema, INSDet generates the test cases by mutating the non-constant fields, including  $R_n$ ,  $R_t$ ,  $P$ ,  $U$ ,  $W$  and  $Imm8$ . This can generate syntactically correct instructions. However, this step is not enough, since it may not generate the values that satisfy the symbolic expression in the ASL code. For instance, one symbolic expression in line 8 of Figure 2b is  $t == 15$ . The random values generated in the first step may not satisfy this expression (all of them are not equal to 15). To this end, we leverage a constraint solver to find the concrete value of the encoding symbol  $R_t$  that satisfies the constraint, i.e., 15. Note that, we only use this to illustrate the basic idea. The concrete value 15 of  $R_t$  likely has been generated in the first step. We take similar actions to solve the constraints for other symbols in line 1 ( $add$ ), 2 ( $index$ ) and 4 ( $wback$ ) of Figure 2c. During this process, we generated 576 instruction streams as test cases in total.

### 3.1.3 Differential testing

We feed each instruction stream into our differential testing engine. The engine generates a corresponding ELF binary for each test case by adding prologue and epilogue instructions. The prologue instructions first set the initial execution context, then the instruction stream will be executed. Finally, the epilogue instructions will dump the result for comparison. We execute the binary on both QEMU and real devices (e.g., RaspberryPi 2B). By comparing the execution result, we confirm that 0xf84f0ddd is an inconsistent instruction stream. Specifically, It will generate a SIGILL signal in a real device while a SIGSEGV signal in QEMU.

We further analyzed the root cause and successfully disclosed a bug in QEMU. According to Figure 2a, the concrete value of the encoding symbol  $R_n$  of the instruction stream 0xf84f0ddd is 1111. As shown in the ASL code (line 1) in Figure 2b, it is an UNDEFINED instruction stream. However, QEMU does not properly check this condition. Figure 3 shows the (patched) function (i.e., `op_store_r1`) in QEMU for decoding the instruction STR (immediate). It continues the decoding process directly from line 12 without any check. We then submit this bug to QEMU developers and the patch is issued (as shown in line 8-10).

## 3.2 Test Case Generator

In theory, for a 32-bit instruction, there exist  $2^{32} = 4,294,967,296$  possible instruction streams, which are not practical for evaluation. In our work, we need to generate a small number of representative test cases that cover most behaviors of an instruction.

Specifically, we first parse the encoding schema to retrieve the encoding symbols and then infer the type for symbols, e.g., a register index or an immediate value. After that, we generate an initialized mutation set by pre-defined rules for each type of the symbol (Table 1 shows the detailed rules). For instance, we generate the maximum, minimum and random values for an immediate value. Then, we develop a symbolic execution engine to solve the constraints in the ASL code for the decoding and execution logic. This step can add more values to the mutation set to satisfy the constraints of the symbols in the ASL program. At last, we remove duplicate values and then generate instruction streams as test cases.

Algorithm 1 shows how we generate the test cases. For each instruction, ARM provides a XML file to describe the instruction. We extract the encoding schemas and the corresponding ASL code for decoding and execution by parsing

```

1  static bool op_store_ri(DisasContext *s, arg_ldst_ri *a, MemOp mop, int mem_idx)
2  {
3      ISSInfo issinfo = make_issinfo(s, a->rt, a->p, a->w) | ISSIsWrite;
4      TCGv_i32 addr, tmp;
5
6      // Rn=1111 is UNDEFINED for Thumb;
7
8      + if (s->thumb && a->rn == 15) {
9      +     return false;
10     + }
11
12     addr = op_addr_ri_pre(s, a);
13
14     tmp = load_reg(s, a->rt);
15     gen_aa32_st_i32(s, args);
16     disas_set_da_iss(s, mop, issinfo);
17     tcg_temp_free_i32(tmp);
18     op_addr_ri_post(s, a, addr, 0);
19     return true;
20 }

```

Figure 3: Original code of QEMU and the patch for function `op_store_ri`, which aims to translate STR instruction

---

**Algorithm 1:** The algorithm to generate test cases.

---

**Input:** The encoding diagram:  $I\_Encode$ ;

The decoding ASL code:  $I\_Decode$ ;

The execution ASL code:  $I\_Execute$

**Output:** The generated test cases:  $T$ ;

```

1  Function Generate( $I\_Encode, I\_Decode, I\_Execute$ ):
2       $Symbols, Constants, Constraints = ParseASL(I\_Encode, I\_Decode, I\_Execute)$ 
3      for  $S$  in  $Symbols$  do
4           $S.MutationSet = InitSet(S)$ 
5      for  $C$  in  $Constants$  do
6           $C.MutationSet = [ConstantValue]$ 
7      for  $C$  in  $Constraints$  do
8           $ValueSet = SolveConstraint(C, Symbols, I\_Decode, I\_Execute)$ 
9          for  $V, S$  in  $ValueSet$  do
10             if  $V$  not in  $S.MutationSet$  then
11                  $S.MutationSet$  add  $V$ 
12      $MutationSets = [S.MutationSet + C.MutationSet]$ 
13      $TestCase = CartesianProduct(MutationSets)$ 
14     return  $T$ 

```

---

the XML file. We first retrieve the encoding symbols ( $Symbols$ ) and constant values ( $Constants$ ) in the encoding schema, as well as  $Constraints$  for the symbolic expression in decoding and execution ASL code (line 2). We then iterate over the  $Symbols$  and generate the  $MutationSet$  for each symbol (line 3-4), which will be introduced in detail in Section 3.2.1. Note this is the initial mutation set for each symbol. For the  $Constants$ , the  $MutationSet$  contains only the fixed value (line 5-6). After that, we solve the constraints to generate new mutation set (i.e.,  $ValueSet$ ) for each symbol (line 7-8), which will be introduced in detail in Section 3.2.2. Then we check whether the solved value for each symbol is in the symbol's  $MutationSet$  (line 9). If not, we append it to the symbols's  $MutationSet$  (line 10-11). After that, we combine them to get the  $MutationSets$  (line 12).

Finally, considering all the possible combinations of the candidates in the  $MutationSet$  for each symbol, we conduct the Cartesian Product on the  $MutationSets$  to get the test cases for this specific instruction encoding (line 13).

### 3.2.1 Initialize Mutation Set

In the phase of initializing mutation set for each symbol, we consider the types of different symbols and aim to cover different values for different types of symbols. In particular, we infer the type based on the symbol name. For instance, a symbol that represents a register index usually has the name  $Rd$ ,  $Rm$ ,  $Rn$ , etc. As for the immediate value, the symbol name used to be  $immn$  where  $n$  represents the length of the value. For example, the symbol  $imm8$  represents a 8-bit immediate value.

Table 1: The rules of initializing the mutation set.

Type of Symbol Name	Mutation Set
Register Index	0 (R0); 1 (R1); 15 (PC); Random index values
Immediate Value in N bits	Maximum value: $2^N - 1$ ; Minimum value: 0; (N-2) Random Value from the enumerated values
Condition	"1110" (Always execute)
Others in 1 bit	"0"; "1"
Others in N bit (N > 1)	N random value from the enumerated values

31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0
1	1	1	1	0	1	0	0	0	D	1	0	Rn		Vd		0		0	0	x	Size	Align	Rm								
Type																															

(a) Encoding diagram of instruction VLD4 in A32 instruction set

```

1  case type of
2    when '0000'
3      inc = 1;
4    when '0001'
5      inc = 2;
6  if size == '11' then UNDEFINED;
7  alignment = if align == '00' then 1 else 4 << UInt(align);
8  ebytes = 1 << UInt(size);
9  elements = 8 DIV ebytes;
10 d = UInt(D:Vd);
11 d2 = d + inc;
12 d3 = d2 + inc;
13 d4 = d3 + inc;
14 n = UInt(Rn);
15 m = UInt(Rm);
16 wback = (m != 15);
17 register_index = (m != 15 && m != 13);
18 if n == 15 || d4 > 31 then UNPREDICTABLE;

```

(b) Decoding code of instruction VLD4 in A32 instruction set

Figure 4: Test case generator example.

Table 1 shows the rules to initialize the mutation set. For a register index, we include the PC register (index 15), R0, R1 and random values in the set. The register R0 and R1 are used to represent the return value for function calls. As for PC, it can explicitly change the execution flow of the program. Thus, the register index in many instruction encodings cannot be 15. We include it in the mutation set to cover such cases. For the immediate value, the maximum and minimum value are the two boundary values that need to be covered. Apart from this, we randomly select (N-2) values, where N represents the bit length of the symbol. Note that enumerating all the values for one symbol is not realistic because immediate values have 24 bits, resulting in  $2^{24} = 16777216$  candidates.

### 3.2.2 Solve Constraints

Symbolic expressions in ASL code represent the different execution paths of the instruction. For instance,  $d4$  in Figure 4 is a symbolic expression ( $d4 = \text{UInt}(D : Vd) + inc + inc + inc$ ) that determines whether the instruction is an UNPREDICTABLE one. To make our test case representative, the generated test cases should cover as many execution paths as possible. To this end, we design and implement a symbolic execution engine for the ASL code. Specifically, we assign symbolic values for encoding symbols. Then we generate the symbolic expression for each variable in the ASL code. After that, we retrieve the constraint of the symbolic expression and find the concrete values of the encoding symbols that satisfy the constraint and its negation, e.g., solve the constraints  $(d4 > 31) == \text{true}$  and  $(d4 > 31) == \text{false}$ .

Figure 4 shows a concrete example. In line 18, there is a symbolic expression  $d4$  and a constraint  $d4 > 31$ . All the related statements (line 3, 5, 10, 11, 12, and 13) are retrieved via backward slicing and highlighted in the green color. To solve this constraint, we conduct backward symbolic execution. Specifically, the symbol  $d4$  is calculated by the expression  $d4 = d3 + inc$  in line 13. Thus, the constraint is converted to  $d3 + inc > 31$ . Given the relationship between  $d3$  and  $d2$  in line 11, and between  $d2$  and  $d1$  in line 11, we further convert it to  $\text{UInt}(D : Vd) + 3 \times inc > 31$ . The expression  $\text{UInt}(D : Vd)$  is converted to  $Vd + 2^4 \times D$  as the symbol  $Vd$  has 4 bits. Thus, we have the constraint  $Vd + 16 \times D + 3 \times inc > 31$ . Symbol  $inc$  is assigned at line 3 and line 5. Thus, the constraint is  $inc == 1 \text{ or } inc == 2$ .

Apart from this, we need to consider the length of each symbol. Since  $D$  is one bit and  $Vd$  has four bits. Their constraints are  $D \geq 0$  and  $D < 2$ ,  $Vd \geq 0$  and  $Vd < 16$ .

We feed all these constraints to the SMT solver. It returns a solution which is a combination of symbol values that satisfy the constraints. One possible solution is that  $Vd$  is 13,  $D$  is 1, and  $inc$  is 2. We then negate the constraint  $d4 > 31$  and repeat the above mentioned process. In this case, the solution is  $Vd$  is 0,  $D$  is 0 and  $inc$  is 1. Thus, the generated *ValueSet* contains three symbols and each symbol has two candidate values. Note  $inc$ 's value depends on  $Type$ 's value. As we will also solve the constraint  $Type == '0000'$  and  $Type == '0001'$ , the final mutation set of  $Type$  must contain the value that can make  $inc$  to be either 1 or 2. Due to the Cartesian Product between each symbol's mutation set, we can always generate the instruction streams that can satisfy the constraint  $d4 > 31$  and its negation.

Note that the path explosion in symbolic execution is not an issue for our purpose since the decoding and execution ASL code has limited constraints, resulting in limited paths. Meanwhile, we model the utility function calls (e.g., `UInt`) so that the symbol will not be propagated into these functions. Our experiment in Section 4.1 shows that we can generate the test cases within 4 minutes.

### 3.2.3 A Demonstration Example

Table 2 describes how we generate all the test cases for instruction `VLD4` in Figure 4. In total, we split the encoding diagram into nine parts including seven symbols and two constant values (None in the column "Symbol Name"). For constant values, the initialized mutation set has one fixed value. For other symbols, we initialize the mutation set, which is described in column "Init Mutation Set", according to algorithm 1. Then we extract the constraints, and find the satisfied values. Column "Related Constraints" lists the constraints for each symbol. After solving the constraints and their negations, new mutation sets for each symbol will be generated. Finally, we have the mutation set for each symbol, which is denoted by column "Final Mutation Set". We conduct the Cartesian Product between the mutation set of each symbol. In total, we generate  $1 \times 2 \times 1 \times 4 \times 6 \times 2 \times 4 \times 3 \times 5 = 5,760$  test cases for this instruction encoding.

Table 2: The generated mutation set for each symbol of instruction `VLD4` in Figure 4

Symbol Name	Bit Length	Start Offset	End Offset	Type	Init Mutation Set	Related Constraint	Set Added by Solving Constraints and Their Negations	Final Mutation Set	Set Size
None	9	23	31	Fixed Value	"111101000"	NA	NA	"111101000"	1
D	1	22	22	Others in 1 bit	"0", "1"	$d4 > 31$	"0", "1"	"0", "1"	2
None	2	20	21	Fixed Value	"10"	NA	NA	"10"	1
Rn	4	16	19	Register Index	"0000", "0001", "0110", "1111"	$n == 15$	"0000", "1111"	"0000", "0001", "0110", "1111"	4
Vd	4	12	15	Others in 4 bit	"0101", "0110", "1001", "1100"	$d4 > 31$	"0000", "1101"	"0000", "0101", "0110", "1001", "1100", "1101"	6
Type	4	8	11	Others in 4 bit	"0000", "0001"	$Type == '0000'$ $Type == '0001'$	"0000", "0001"	"0000", "0001"	2
Size	2	6	7	Others in 2 bit	"01", "10"	$Size == '11'$	"00", "11"	"00", "01", "10", "11"	4
Align	2	4	5	Others in 2 bit	"00", "11"	$Align == '00'$	"00", "01"	"00", "01", "11"	3
Rm	4	0	3	Register Index	"0000", "0001", "0111", "1111"	$m != 15$ $m != 13$	"0000", "1101", "1111"	"0000", "0001", "0111", "1101", "1111"	5

## 3.3 Differential Testing Engine

### 3.3.1 Model the CPU

The differential testing engine receives the generated test cases (instruction streams), and detects inconsistent ones. Formally, given one instruction stream  $I$ , we denote the state before the execution of  $I$  as the initial state  $CPU_I$  and the state after the execution of  $I$  as the final state  $CPU_F$ . We denote the CPU  $T$ 's initial state  $CPU_I(T)$  with the tuple  $\langle PC_T, Reg_T, Mem_T, Sta_T \rangle$ .  $PC$  denotes the program counter, which points to the next instruction that will be executed.  $Reg$  denotes the registers used by processors while  $Mem$  denotes the memory space that the tested instruction  $I$  may write into. Note we do not consider the whole memory space due to two reasons. First, comparing the whole memory space is time- and resource-consuming. Second, memory addresses like stack address are different each time due to specific memory management strategies (e.g., Address Space Layout Randomization).  $Sta$  denotes the status register, which is *APSR* in ARM architecture.

We denote the CPU  $T$ 's final state  $CPU_F(T)$  with the tuple  $[PC_T, Reg_T, Mem_T, Sta_T, Sig_T]$ . Inside  $CPU_F(T)$ , all the other attributes except  $Sig$  have the same meanings as they are inside  $CPU_I(T)$ .  $Sig$  denotes the signal that the instruction stream  $I$  may trigger. If no signal is triggered, the value of  $Sig$  is 0.

Given the CPU emulator  $E$ , the real device  $R$ , our differential testing engine guarantees that  $E$ 's initial state  $CPU_I(E)$  is equal to the  $R$ 's initial state  $CPU_I(R)$ .  $CPU_I(E) = CPU_I(R)$  iff:



```

1  int main() {
2      register_signals(sig_handler);
3      set_init_state();
4      __asm__("nop");
5      dump_final_state();
6      exit();
7  }
8
9  void sig_handler() {
10     dump_final_state();
11     exit();
12 }

```

Figure 5: The pseudo code for rendering different test cases.

$$\forall \phi \in \langle PC, Reg, Mem, Sta \rangle: \phi_E = \phi_R$$

After the execution of  $I$ ,  $I$  is treated as an inconsistent instruction stream if the final state  $CPU_F(E)$  is not equal to the  $R$ 's final state  $CPU_F(R)$ . More formally,  $CPU_F(E) \neq CPU_F(R)$  iff:

$$\begin{cases} \exists \phi \in [PC, Reg, Mem, Sta] : \phi_E \neq \phi_R & Sig_E = Sig_R = 0 \\ True & Sig_E \neq Sig_R \\ \exists \phi \in [PC, Reg, Sta] : \phi_E \neq \phi_R & Sig_E = Sig_R \neq 0 \end{cases}$$

This is because if the instruction stream  $I$  triggers a signal on both CPU emulator  $E$  and a real device  $R$  with a different signal number,  $CPU_F(E)$  is not equal to  $CPU_F(R)$ . If the signal number is the same, we need to compare the  $PC$ ,  $Reg$ , and  $Sta$  as the instruction stream is not executed normally (otherwise no signal will be raised). When the  $I$  does not trigger a signal on a CPU emulator  $E$  or a real device  $R$  ( $Sig_E = Sig_R = 0$ ), we will compare  $Mem$ .

Note we do not change any internal logic of the CPU emulator. Thus, our differential testing engine can be principally applied to different emulators, if necessary.

### 3.3.2 Our Strategy

To conduct the differential testing, we build a template binary  $B$ . Given one instruction  $I$ , we will generate a new binary  $B_I$  by inserting prologue and epilogue instructions. Figure 5 shows pseudo code. We first register the signal handlers (line 2) to capture different signals. This is because the test instruction may trigger exceptions such as illegal encoding, memory error, etc.

To make the initial state consistent (line 3), we set the value of general purpose registers to zero except the frame register (R11), the stack register (R13), the link register (R14), and the PC. This is because these registers have specific functionalities. For instance, link register is used to save the return address (a non-zero value) while PC can influence the execution flow.

After setting up the initial state, a test case (an instruction stream) will be executed. We set the instruction `nop` with inline assembly (line 4). Given one test instruction stream, we statically rewrite the binary and change the `nop` instruction to the test instruction stream to generate a new binary. After the execution, we need to dump the CPU state (line 5 and line 10) so that we can compare the execution result. Note the test instruction stream could trigger a signal, thus, we also need to dump the result in signal handlers.

For register values, we push them on the stack and then write them into a file. For the memory, we check the test instruction with Capstone [4], analyze the instruction to see whether it will write a value into a memory location. If so, we load the memory address, and push it on the stack for later inspection. Finally, we compare the result collected from the emulator and a real device. If the instruction stream results in a different CPU final state, ( $CPU_F(E) \neq CPU_F(R)$ ), it will be treated as an inconsistent instruction stream.

## 3.4 Implementation Details

We have implemented a prototype system INSDet using Python, C and ARM assembly. In particular, we implement the test case generator in Python. We parse the ASL code, extract the lexical and syntactic information with regular expressions. We use Z3 [8] as the SMT solver to solve the constraints. The differential testing engine is implemented in C and assembly code with some glue scripts in Python. Specifically, the initial state setup and the execution result

Table 3: The statistics of the generated instructions. "GIS" denotes the number of generated instruction streams. "VIS" denotes the number of valid instruction streams, which means they match the instruction encodings. "VISR" denotes the percentage of dividing "VIS" by "GIS". "AE" denotes the number of all instruction encodings. "CE" denotes the number of covered instruction encodings by the generated instruction streams. "CER" denotes the percentage of dividing "CE" by "AE". "AI" denotes the number of all instructions. "CI" denotes the number of covered instruction encodings by the generated instruction streams. CIR denotes the percentage of dividing "CI" by "AI". Note one instruction may have different instruction encodings for different instruction set. The total instruction for instruction set A32, T32, and T16 is 489.

Instruction Set	Time (s)	GIS	VIS	VISR	AE	CE	CER	AI	CI	CIR	Solved Constraints
A64	70.51	1,094,700	1,094,700	100	839	839	100	581	581	100	3,436
A32	75.05	870,221	870,221	100	550	550	100	481	481	100	4,718
T32	74.58	808,770	808,770	100	531	531	100	451	451	100	4,425
T16	2.32	958	958	100	78	78	100	68	68	100	122
Overall	222.46	2,774,649	2,774,649	100	1,998	1,998	100	1,070	1,070	100	12,701

dumping is implemented with inline assembly. In total, INSDet contains 5,074 lines of Python code, 220 lines of C code, and 200 lines of assembly code.

## 4 Evaluation

In this section, we evaluate INSDet by answering the following three research questions.

- **RQ1:** Is INSDet able to generate sufficient test cases?
- **RQ2:** Is INSDet able to detect inconsistent instructions? What are the root causes of these inconsistent instructions?
- **RQ3:** What are the possible usage scenarios of inconsistent instructions?

### 4.1 Sufficiency of Test Case Generator (RQ1)

We generate the test cases according to ARMv8-A manual, which introduces ASL. Specifically, the manual includes four different instruction sets. In AArch 64 mode, A64 instruction set is supported. For the AArch 32 mode, it consists of three different instruction sets. They are ARM32 with 32-bit instruction length (A32), Thumb-2 with instruction length of mixed 16-bits and 32-bits (T32), and Thumb-1 with 16-bit instruction length (T16). They are also supported by previous ARM architectures (e.g., ARMv6, ARMv7). To locate the inconsistent instructions in different ARM architectures, we generate the test cases for all the instruction sets.

The generated test case is sufficient. Table 3 shows the statistics of the generated instructions. For A64, we generate around 1 million instructions, which cover all the 839 instruction encodings in 581 instructions. In the decoding and executing ASL code, we solved 3,436 constraints that are related to encoding symbols.

We noticed that both A32 and T32 have around 5 hundred instruction encodings and more than 800 thousands instructions are generated. For T16, the generated instruction streams are rather small due to the small number of instruction encoding schemes and limited instruction length. Overall, all the generated instruction streams are valid instruction streams, which means they meet the encoding schema of one instruction encoding. Meanwhile, all the instruction encodings and instructions are covered. Furthermore, more than 12 thousand constraints, which are related to encoding symbols, are solved, indicating the multiple behaviors of the instructions are explored.

To demonstrate the effectiveness of the test case generation algorithm, we randomly generate some instruction streams. To make the comparison fair, we generate the same number of test cases for each instruction set. Table 4 shows the result. We repeat the randomly generated process for 10 times. Then we check whether the generated instructions are valid instruction streams or not. If they are the valid instruction streams, we calculate how many instruction encodings, how many instructions, and how many constraints are covered by these instruction streams. We noticed that only 37.3% generated instruction streams are valid instruction streams, which means all the other instructions are illegal instructions and they are not effective to test the potential different behaviors between real devices and CPU emulators. Among the valid instruction streams, it can only cover 54.5% instruction encodings and 51.4% instructions. Nearly a half of instructions can not be covered with the randomly generated instruction streams. Specifically, many of the T32 instructions cannot be covered with randomly generated instructions, which means many of these instructions have fixed value. As for the coverage of constraints, 62.6% constraints are covered while the left 37.4% constraints can not be explored, resulting in a relatively low behavior space.

Table 4: The statistics of the generated instructions in random. "IS" denotes instruction set. "Ave" denotes average value. For each instruction set and the overall result, "R" inside column "VIS", "CE", "CI", and "CC" denotes the percentage of dividing "Ave" inside the same column by "VIS", "CE", "CI", and "Solved Constraints" in Table 3, respectively. "GIS", "VIS", "CE", "CI" denotes the same meaning illustrated in Table 3. "CC" denotes covered constraints.

IS	GIS	VIS		CE		CI		CC	
		Ave	R	Ave	R	Ave	R	Ave	R
A64	1,094,700	421,645	38.5	265	31.6	178	30.6	934	27.2
A32	870,221	578,845	66.5	415	75.5	361	75.1	3,725	79.0
T32	808,770	34,598	4.2	351	66.1	283	62.7	3,203	72.3
T16	958	796	83.0	57	73.1	49	72.1	84	68.9
Overall	2,774,649	1,035,884	37.3	1,088	54.5	550	51.4	7,946	62.6

**Answer to RQ1:** INSDet can generate sufficient test cases, which are all valid instruction streams and can cover all instruction encodings and instructions. On the contrary, Only 37.3% of the same number of randomly generated instruction streams are valid instruction streams. Furthermore, 45.5% instruction encodings, 48.6% instructions, and 37.4% constraints cannot be explored by these randomly generated instructions.

## 4.2 Differential Testing Results and Root Causes (RQ2)

With the generated test cases in four different instruction set. We feed them into our differential testing engine to locate the inconsistent instructions. Table 5 shows the result.

**Experiment Setup** In total, we conduct the differential testing between QEMU (version 5.1.0) and four different ARM architecture versions (i.e., ARMv5, ARMv6, ARMv7-a, ARMv8-a). For each ARM architecture version, we select one real device. Specifically, we select devices OLinuXino iMX233 for ARMv5, RaspberryPi Zero for ARMv6, RaspberryPi 2B for ARMv7-a, and Hikey 970 for ARMv8-a. To make the differential testing fair, we use the same or similar CPU model between the emulator and real device. Note the CPU model for Hikey 970 is A73 while the most advanced CPU model supported by QEMU is A72, which is selected. However, they two both share the same instruction set (i.e., ARMv8-a). For ARMv5, only A32 instruction set is supported. Since QEMU does not support Thumb2 for ARM1176 of ARMv6, we only test the A32 instruction set on ARMv6. For ARMv7, all the instruction set is supported except A64. Since the T16 instruction has a rather small number of set. We combine the T16 and T32 in the testing process. For ARMv8-A, only A64 instruction set is supported in user-level programs. The "Generated Instruction Streams", "All Instruction Encodings", and "All Instructions" in Table 5 are from the "GIS", "AE", and "AI" in Table 3, respectively. In total, it takes around 2700 seconds of CPU time for QEMU, which is run on the Intel i7-9700 CPU. For the real devices, the CPU time cost ranges from 5276 seconds to 46238 seconds (around 13 hours), depending on the specific devices. Thanks to the representative test cases, the differential testing for all the test cases can be finished within acceptable time.

**Testing Result** Among all the test cases, some of them may read from or write into SP and FP registers. SP and FP are used to store function parameters and local variables. They are influenced by the memory management strategy (e.g., Address Space Layout Randomization) of the whole system and are different for each run. Meanwhile, some instruction streams may change these two registers and result in crash of the whole test binary. Thus, instruction streams that read from or write into these two registers are filtered. Apart from this, some instruction streams are branch instructions. These instructions may execute normally on both emulator and real devices. Then inconsistent behaviors may occur due to the other instructions as the branch instructions change the execution flow of the binary. We also filter these instructions if they execute normally. For the left instruction streams, we noticed around 90% instruction streams (from 88.9% to 94.3% for each architecture version) are consistent for four different ARM architectures. However, there are still thousands of inconsistent instruction streams. In particular, 155,642 inconsistent instruction streams are found, owning to 5.6% of the whole test cases. Note one instruction may be tested in different architectures (e.g., A32 instruction set in ARMv5, ARMv6, and ARMv7), the number in column "Overall" is the union of the other columns, which means the instruction stream can cause inconsistent behavior in at least one architecture. Furthermore, these inconsistent instruction streams cover 600 different instruction encodings and 511 instructions, owning 30% and 47.8% of the all instruction encodings and instructions, respectively.

**Inconsistent Behavior** We further analyze the inconsistent instruction streams and categorize them according to our modeled CPU. We noticed nearly half of the inconsistent instruction streams (i.e., 47.3%) raise different signal numbers during the execution. Meanwhile, 42.1% inconsistent instruction streams would trigger signals in real devices while they can be executed normally in QEMU. The percentage of inconsistent instruction streams drops to 4.6%

Table 5: The results of differential testing. "CPU Time" denotes the sum of the time used by CPU for all the test cases, which is in seconds. We do not count the overall CPU time for real devices as different devices have different CPUs. "Inst": Instructions; "Enc": Encodings; "Cons\_Unpre": Constraint UNPREDICTABLE; "Annotation\_Def": Defined in Annotation;  $\langle X | Y \rangle$ : For each instruction set and overall result, X denotes the number of the attribute indicated by the row name while Y denotes the percentage of dividing X by Z. For data in "Testing Result", Z stands for the row "Generated Instruction Streams", "All Instruction Encodings", and "All Instructions". For data in "Inconsistent Behaviors" and "Root Cause", Z stands for "Inconsistent Instruction Streams".  $[M | N]$ : For each instruction set and overall result, M denotes the number of instruction encodings while N denotes how many instructions M belongs to. We do not calculate the percentage of instruction encodings or instructions for "Inconsistent Behaviors" and "Root Cause" as one encoding or instruction can have more than one inconsistent behaviors, resulting from different root causes.

Architecture	ARMv5	ARMv6	ARMv7-a		ARMv8-a	Overall
Experiment Setup						
QEMU Binary	qemu-arm	qemu-arm	qemu-arm		qemu-aarch64	-
QEMU Model	ARM926	ARM1176	Cortex-A7		Cortex-A72	-
Device Name	OLinuXino IMX233	RaspberryPi Zero	RaspberryPi 2B		Hikey 970	-
Device Model	ARM926	ARM1176	Cortex-A7		Cortex-A73	-
Instruction Set	A32	A32	T32&T16		A64	-
Generated Instruction Streams	870,221	870,221	870,221	809,728	1,094,700	2,774,649
All Instruction Encodings	550	550	550	609	839	1,998
All Instructions	481	481	481	462	581	1,070
CPU Time (QEMU)	530.5	540.6	538.2	462.1	625.9	2697.3
CPU Time (Device)	46238.0	6901.7	6194.2	5276.0	9145.0	-
Testing Result	The percentage is based on the number of generated instructions and all encodings					
Read/Write to SP/FP	$\langle 37,002   4.3\% \rangle$	$\langle 36,879   4.2\% \rangle$	$\langle 37,002   4.3\% \rangle$	$\langle 24,582   3.0\% \rangle$	$\langle 45,985   4.2\% \rangle$	$\langle 107,569   3.9\% \rangle$
Branch	$\langle 3,557   0.4\% \rangle$	$\langle 5,431   0.6\% \rangle$	$\langle 1,821   0.2\% \rangle$	$\langle 1,048   0.1\% \rangle$	$\langle 210   0.0\% \rangle$	$\langle 7,017   0.4\% \rangle$
Consistent Instruction Streams	$\langle 794,418   91.3\% \rangle$	$\langle 818,744   94.1\% \rangle$	$\langle 773,906   88.9\% \rangle$	$\langle 738,369   91.2\% \rangle$	$\langle 1,031,901   94.3\% \rangle$	$\langle 2,598,695   93.7\% \rangle$
Inconsistent Instruction Streams	$\langle 35,244   4.1\% \rangle$	$\langle 9,167   1.1\% \rangle$	$\langle 57,492   6.6\% \rangle$	$\langle 45,729   5.6\% \rangle$	$\langle 16,604   1.5\% \rangle$	$\langle 155,642   5.6\% \rangle$
Inconsistent Instruction Encodings	$\langle 108   19.6\% \rangle$	$\langle 86   15.6\% \rangle$	$\langle 270   49.0\% \rangle$	$\langle 266   43.6\% \rangle$	$\langle 15   1.8\% \rangle$	$\langle 600   30\% \rangle$
Inconsistent Instructions	$\langle 98   20.4\% \rangle$	$\langle 84   17.5\% \rangle$	$\langle 229   47.6\% \rangle$	$\langle 223   48.3\% \rangle$	$\langle 13   2.2\% \rangle$	$\langle 511   47.8\% \rangle$
Inconsistent Behaviors	The percentage is based on the number of inconsistent instructions					
$Sig_E \neq Sig_R \neq 0 [I.S   I.S.R]$	$\langle 10,748   30.6\% \rangle$	$\langle 6,853   74.8\% \rangle$	$\langle 35,628   62.0\% \rangle$	$\langle 24,094   52.7\% \rangle$	$\langle 667   4.0\% \rangle$	$\langle 73,619   47.3\% \rangle$
$Sig_E \neq Sig_R \neq 0 [Enc   Inst]$	$\langle 57   51 \rangle$	$\langle 63   61 \rangle$	$\langle 209   186 \rangle$	$\langle 212   181 \rangle$	$\langle 7   5 \rangle$	$\langle 458   401 \rangle$
$Sig_E \neq Sig_R = 0 [I.S   I.S.R]$	$\langle 4,559   13.0\% \rangle$	$\langle 1,433   15.6\% \rangle$	$\langle 551   1.0\% \rangle$	$\langle 836   1.8\% \rangle$	$\langle 0   0.0\% \rangle$	$\langle 7,197   4.6\% \rangle$
$Sig_E \neq Sig_R = 0 [Enc   Inst]$	$\langle 34   33 \rangle$	$\langle 47   46 \rangle$	$\langle 18   18 \rangle$	$\langle 9   9 \rangle$	$\langle 0   0 \rangle$	$\langle 72   70 \rangle$
$Sig_R \neq Sig_E = 0 [I.S   I.S.R]$	$\langle 13,827   39.3\% \rangle$	$\langle 3   0.0\% \rangle$	$\langle 21,135   36.8\% \rangle$	$\langle 19,915   43.6\% \rangle$	$\langle 11,219   67.6\% \rangle$	$\langle 65,599   42.1\% \rangle$
$Sig_R \neq Sig_E = 0 [Enc   Inst]$	$\langle 42   35 \rangle$	$\langle 1   1 \rangle$	$\langle 193   168 \rangle$	$\langle 203   177 \rangle$	$\langle 7   7 \rangle$	$\langle 420   367 \rangle$
$Sig_R = Sig_E \neq 0 [I.S   I.S.R]$	$\langle 0   0.0\% \rangle$	$\langle 67   0.7\% \rangle$	$\langle 174   0.3\% \rangle$	$\langle 873   1.9\% \rangle$	$\langle 4,716   28.4\% \rangle$	$\langle 5,763   3.7\% \rangle$
$Sig_R = Sig_E \neq 0 [Enc   Inst]$	$\langle 0   0 \rangle$	$\langle 6   6 \rangle$	$\langle 18   18 \rangle$	$\langle 18   15 \rangle$	$\langle 3   3 \rangle$	$\langle 39   36 \rangle$
$Sig_R = Sig_E = 0 [I.S   I.S.R]$	$\langle 2,389   6.8\% \rangle$	$\langle 318   3.5\% \rangle$	$\langle 3   0.0\% \rangle$	$\langle 8   0.0\% \rangle$	$\langle 0   0.0\% \rangle$	$\langle 2,410   1.5\% \rangle$
$Sig_R = Sig_E = 0 [Enc   Inst]$	$\langle 25   25 \rangle$	$\langle 6   6 \rangle$	$\langle 1   1 \rangle$	$\langle 1   1 \rangle$	$\langle 0   0 \rangle$	$\langle 27   27 \rangle$
Others $[I.S   I.S.R]$	$\langle 3,721   10.6\% \rangle$	$\langle 493   5.4\% \rangle$	$\langle 1   0.0\% \rangle$	$\langle 3   0.0\% \rangle$	$\langle 2   0.0\% \rangle$	$\langle 4,218   2.8\% \rangle$
Others $[Enc   Inst]$	$\langle 8   8 \rangle$	$\langle 2   2 \rangle$	$\langle 1   1 \rangle$	$\langle 3   2 \rangle$	$\langle 2   2 \rangle$	$\langle 14   13 \rangle$
Root Cause	The percentage is based on the number of inconsistent instructions					
Bugs of QEMU $[I.S   I.S.R]$	$\langle 1   0.0\% \rangle$	$\langle 1   0.0\% \rangle$	$\langle 1   0.0\% \rangle$	$\langle 583   1.3\% \rangle$	$\langle 2   0.0\% \rangle$	$\langle 586   0.4\% \rangle$
Bugs of QEMU $[Enc   Inst]$	$\langle 1   1 \rangle$	$\langle 1   1 \rangle$	$\langle 1   1 \rangle$	$\langle 10   7 \rangle$	$\langle 2   2 \rangle$	$\langle 13   10 \rangle$
UNPREDICTABLE $[I.S   I.S.R]$	$\langle 35,147   99.7\% \rangle$	$\langle 8,913   97.2\% \rangle$	$\langle 57,428   99.9\% \rangle$	$\langle 44,869   98.1\% \rangle$	$\langle 2   0.0\% \rangle$	$\langle 137,828   88.6\% \rangle$
UNPREDICTABLE $[Enc   Inst]$	$\langle 106   96 \rangle$	$\langle 77   75 \rangle$	$\langle 265   224 \rangle$	$\langle 261   218 \rangle$	$\langle 0   0 \rangle$	$\langle 570   483 \rangle$
Cons.UNPRE $[I.S   I.S.R]$	$\langle 0   0.0\% \rangle$	$\langle 0   0.0\% \rangle$	$\langle 0   0.0\% \rangle$	$\langle 0   0.0\% \rangle$	$\langle 16,602   100.0\% \rangle$	$\langle 16,602   10.7\% \rangle$
Cons.UNPRE $[Enc   Inst]$	$\langle 0   0 \rangle$	$\langle 0   0 \rangle$	$\langle 0   0 \rangle$	$\langle 0   0 \rangle$	$\langle 15   13 \rangle$	$\langle 15   13 \rangle$
Annotation_Def $[I.S   I.S.R]$	$\langle 0   0.0\% \rangle$	$\langle 253   2.8\% \rangle$	$\langle 63   0.1\% \rangle$	$\langle 277   0.6\% \rangle$	$\langle 0   0.0\% \rangle$	$\langle 530   0.3\% \rangle$
Annotation_Def $[Enc   Inst]$	$\langle 0   0 \rangle$	$\langle 8   8 \rangle$	$\langle 4   4 \rangle$	$\langle 4   4 \rangle$	$\langle 0   0 \rangle$	$\langle 12   12 \rangle$
Others $[I.S   I.S.R]$	$\langle 96   0.3\% \rangle$	$\langle 0   0.0\% \rangle$	$\langle 0   0.0\% \rangle$	$\langle 0   0.0\% \rangle$	$\langle 0   0.0\% \rangle$	$\langle 96   0.0\% \rangle$
Others $[Enc   Inst]$	$\langle 1   1 \rangle$	$\langle 0   0 \rangle$	$\langle 0   0 \rangle$	$\langle 0   0 \rangle$	$\langle 0   0 \rangle$	$\langle 1   1 \rangle$

if they can trigger signals on QEMU but cannot on real devices, which demonstrate that QEMU is more tolerant compared with real devices for many instructions. Overall, the above mentioned three cases ( $Sig_E \neq Sig_R$ ) take 94.0% (47.3 + 4.6 + 42.1) of all the inconsistent instruction streams. A small number of instruction streams may (i.e., 3.7%) or may not (i.e., 1.5%) trigger the same signals but have different register or memory values. They mainly because of the UNPREDICTABLE conditions. The left 2.8% inconsistent instruction streams are due to the other problems. These instruction can make the emulator or real devices crash or stuck.

Note instruction streams generated from one instruction encoding can result in different inconsistent behaviors due to the detail decoding and executing logic. According to our experiments, the inconsistent behaviors with more inconsistent instruction streams usually cover more instruction encodings and instructions. For example, 458 different instruction encodings from 401 instructions are covered by the 73,619 inconsistent instruction streams that trigger different signals.

**Root Cause** Based on the above mentioned inconsistent behavior, we explore the root cause of the inconsistent instructions. First, there are implementation bugs of QEMU. We discovered 4 bugs in total, which come from 586 inconsistent instruction streams and 13 instruction encodings. Some of the bugs are related to very common instructions. For example, many load and store instructions (e.g., LDRD, STRD, LDM, STM, etc.) in A32 instruction set should check the alignments mandatory while QEMU does not check. Instructions like STR, BLX may be undefined instructions in specific cases, which should raise SIGILL signal, even they may meet the corresponding encoding schema. However,

```

1  boolean AArch32.ExclusiveMonitorsPass(bits(32) address, integer size)
2  // It is IMPLEMENTATION DEFINED whether the
3  // detection of memory aborts happens before or
4  // after the check on the local Exclusive Monitor.
5  // As a result, a failure of the local monitor can
6  // occur on some implementations even if the
7  // memory access would give an memory abort.
8  ...
9  return

```

Figure 6: Two different implementations are defined in the annotation of function ExclusiveMonitorsPass, which is called by many instructions’ executing code

Table 6: The statistics on detecting emulators

Mobile Type	CPU	A64	A32	T32 & T16
Samsung S8	SnapDragon 835	✓	✓	✓
Huawei Mate20	Kirin 980	✓	✓	✓
IQOO Neo5	SnapDragon 870	✓	✓	✓
Huawei P40	Kirin 990	✓	✓	✓
Huawei Mate40 Pro	Kirin 9000	✓	✓	✓
Honor 9	Kirin 960	✓	✓	✓
Honor 20	Kirin 710	✓	✓	✓
Blackberry Key2	SnapDragon 660	✓	✓	✓
Google Pixel	SnapDragon 821	✓	✓	✓
Samsung Zflip	SnapDragon 855	✓	✓	✓
Google Pixel3	SnapDragon 845	✓	✓	✓

QEMU does not follow the specification. We also noticed one instruction (i.e., WFI) that can make QEMU crash. WFI denotes waiting for interrupt and is usually used in system-mode emulation. However, ARM manual specifies that it can also be used in user-space. QEMU does not handle this instruction well and an abort will be generated during user-space emulation. After our report, all of these bugs are confirmed by developers and are in patching process. This also demonstrates the capability of INSDet in discovering the bugs of the emulator implementation.

Apart from the bugs, many other inconsistent instruction streams are due to the undefined implementation in the ARM manual. There are three different kinds of undefined implementation. The first one is UNPREDICTABLE, which is introduced in Section 3.1. UNPREDICTABLE leaves open implementation decision for emulators and processors, which can result in inconsistent instructions. UNPREDICTABLE is the major reason and it takes 88.6% for all the inconsistent instruction streams. The second is Constraint UNPREDICTABLE (“Cons\_UNPRE” in Table 5). Constraint UNPREDICTABLE provides candidate implementation strategies and the developer or vendor can choose from one of them, which only exists in A64 instructions. The last one is defined in the annotation part of the ASL code (“Annotation\_Def” in Table 5). Figure 6 shows an example. In the function ExclusiveMonitorsPass, which is called by the executing code of instruction STREXH, there is an annotation for the implementation. Note the check on the *local Exclusive Monitor* would update the value of a register. Thus, if the detection of memory aborts happens before the check, the value of the register would not be updated while the detection happens after the check can update the value, resulting in different register value.

Some specific instructions (i.e., BKPT) can trigger SIGTRAP signal in QEMU while make the real devices stuck. As we cannot verify the detail implementation logic of real devices without the design specification, we left it to be others.

**Answer to RQ2:** INSDet can detect inconsistent instructions. In total, 155,642 inconsistent instruction streams are found, which covers 30% (i.e., 600/1998) instruction encodings and 47.8% instructions (i.e., 511/1070). The implementation bugs of QEMU and the undefined implementation in ARM manual are the major root causes. Four bugs are discovered and confirmed by QEMU developers. These bugs influence 13 instruction encodings including commonly used instructions (e.g., STR, BLX).

### 4.3 Applications of Inconsistent Instructions (RQ3)

According to the evaluation result in Section 4.2, there are many inconsistent instructions (i.e., 47.8%) in ARM architectures. These instructions can be used to detect the existence of emulators. Furthermore, detecting emulator can prevent the binary from being analyzed or fuzzed, which is known as anti-emulation and anti-fuzzing technique.

```

1 void sig_handler(int signum) {
2     record_execution_result(i++);
3     siglongjmp(sig_env, i);
4 }
5
6 Bool JNI_Function_Is_In_Emulator() {
7     register_signals(sig_handler);
8     i = sigsetjmp(sig_env, 0);
9     switch (i){
10         case 1:
11             execute(inconsistent_instruction_n);
12             record_execution_result(i++);
13             longjmp(sig_env, i++);
14         case 2:
15             ...
16         case n:
17             }
18     return compare_result();
19 }

```

Figure 7: Pseudo code of the native code for detecting the emulator.

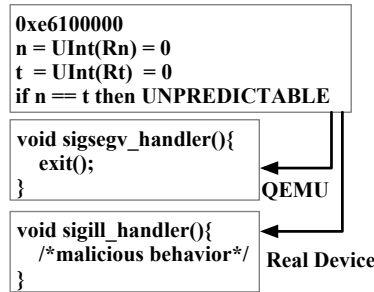


Figure 8: Inconsistent instruction can prevent the malicious behavior being detected by emulators

#### 4.3.1 Emulator Detection

The inconsistent instructions can be used to detect emulators. We use the inconsistent instructions for emulator detection. Considering the popularity of Android systems, we target Android applications. Specifically, we build a native library by using the inconsistent instructions. For each encoding of one inconsistent behavior, we select one instruction stream.

Figure 7 shows the pseudo code of the library. Function *JNI\_Function\_Is\_In\_Emulator* (line 6) returns True if the emulator is detected. Inside the function, we register signal handlers for different signals (line 7). After the execution of each instruction stream, we will record the execution result either in the signal handler (line 2) or after the execution (line 12), which depends on the execution result. Then we use the function *longjmp* (line 13) or *siglongjmp* (line 3) to jump back to the place where calling *sigsetjmp* (line 8). Note the return value of *sigsetjmp* depends on the parameter of the *longjmp* or *siglongjmp*. Thus, the *switch* would check the value of *i*, which is the return value of *sigsetjmp*, repeatedly. As *i* would increase one after the execution of one instruction stream, we can execute hundreds of instruction streams in one function by adding corresponding *case* conditions. Each instruction stream can make an equal contribution to the final decision on whether the current execution environment is in real devices or emulators. Finally, if more instruction streams decide the application are running inside an emulator, the *compare\_result()* will return True and vice versa.

We automatically generate the test library with template code and build three Android apps for different instruction set (one for A64, one for A32, and one for T32 & T16). We run the applications on 12 different mobiles from 6 different vendors. These mobiles also use different CPUs, which decides how an instruction stream should be executed. Meanwhile, we run the applications in the Android emulator provided by Android studio (Version 4.1.2). We compare the result running in each real mobiles with the one in emulators. If the function *JNI\_Function\_Is\_In\_Emulator* returns True in emulator and returns False in real mobiles. We consider it will successfully detect the emulator. Table 6 shows the evaluation result, by testing the three Android apps (one for A64, one for A32, and one for T32 & T16) in 12 mobiles, all the mobile apps can detect the existence of emulator and real mobiles successfully.

```

1 0x10000: e51b3008 LDR r3,[fp,#-8]
2 0x10004: e1a03000 MOV r3,r0
3 0x10008: e7cf0e9f BFC r0, #0xf, #1
4 // BFC instruction is to clear specific bits
5 // e7cf0e9f is an UNPREDICTABLE encoding
6 // e7cf0e9f is executed normally in real device
7 // e7cf0e9f triggers SIGILL signal on QEMU
8 0x1000c: e1a00003 MOV r0,r3
9 0x10010: e50b3008 STR r3,[fp,#-8]

```

Figure 9: Instrumented instruction streams for anti-fuzzing.

### 4.3.2 Anti-Emulation

Anti-emulation technique is important. On the attacker’s side, it can be proposed to increase the bar for analyzing the malware so that the defense mechanism can be developed slower. On the defender’s side, commercial software needs to protect the core functionality and algorithms from being analyzed. Thus, it is widely used in the wild [43].

The inconsistent instructions can be used to conduct anti-emulation and can prevent the malware’s malicious behavior being analyzed. We demonstrate how the inconsistent instruction can be used to hide the malicious behavior.

We use one of the state-of-the-art dynamic analysis platforms (i.e., PANDA [5]) to demonstrate the usage. PANDA is built upon QEMU and supports taint analysis, record and replay, operating system introspection, and so on. We port one of the open source rootkits (i.e., Suterusu [6]) to Debian 7.3. We register two different signal handlers for SIGILL and SIGSEGV, respectively. Then we instrument one instruction stream (i.e., 0xe6100000). This is a LDR instruction encoding in ARM instruction set. According to the encoding schema, n equals to t and both these two symbols’ values are zero. The ASL code of decoding would check whether n equals to t. If so, it should be the UNPREDICTABLE behavior. Real devices think this is an illegal instruction stream and will raise the SIGILL signal while QEMU tries to execute the instruction stream. Then SIGSEGV will be raised as the address pointed by R0 cannot be accessed. In this case, the malicious behavior will only be triggered in real devices. Meanwhile, when we use the PANDA to analyze the malware, no malicious behavior will be monitored and the program will exit inside the *sigsegv\_handler*.

### 4.3.3 Anti-Fuzz

Fuzzing is widely used to explore the zero-day vulnerabilities. To help the released binaries from being fuzzed by attackers, researchers utilize anti-fuzzing techniques [24, 17]. Considering that many new binary fuzzing frameworks are based on QEMU, the inconsistent instructions can be used by developers as a mitigation approach towards fuzzing technique.

We demonstrate how the inconsistent instructions can be used to conduct anti-fuzzing tasks with a relatively low overhead and high decreased coverage ratio.

Figure 9 shows a snippet of assembly code instrumented into the release binary. In address 0x10008, the instruction BFC is used to clear bits for register R0. Note we move the value of R0 to R3 before the instruction BFC and return it back after the execution of BFC. This can guarantee the instrumented instructions will not affect the execution of the binary on the real device. The instruction stream 0xe7cf0e9f results in an UNPREDICTABLE condition. It can be executed normally in real devices while triggering a signal on QEMU.

We developed a GCC plugin to instrument the above mentioned inconsistent instruction streams at each function entry and apply this plugin on three popular used libraries (i.e., libtiff, libpng, and libjpeg) during the compilation process to generate released binaries.

Table 7 shows the space and runtime overhead of the instrumented binary compared with the normal (non-instrumented) one. The space overhead is measured by comparing the binary size. For runtime overhead, we measure it by running test suites on both binaries and comparing the cost of time. We noticed that the instrumented binary imposes negligible space and runtime overhead to the binary. The average space overhead for the protected binary is around 4%, and the runtime overhead is less than 1%.

We then measure the functionality of anti-fuzzing. We fuzz the instrumented binaries and the normal ones with AFL-QEMU (version 2.56b) for 24 hours. The seed corpus is the test suite used for each library in Table 7. We collect the coverage information for the instrumented and the normal ones. Figure 10 shows the results. It is easy to see that the coverage for instrumented binaries cannot increase (because QEMU fails to execute binaries correctly), while the normal ones will increase with the fuzzing time.

Table 7: Overhead information of anti-fuzzing.

Library <sup>1</sup>	Test Suite <sup>2</sup>	Space Overhead	Runtime Overhead
libpng (readpng)	built-in (254)	4.0% (+7KB)	0.52%
libjpeg (djpeg)	GIT <sup>3</sup> (97)	4.3% (+8KB)	0.61%
libtiff (tiffinfo)	built-in (61)	2.2% (+8KB)	0.59%
Overall		3.5%	0.57%

<sup>1</sup> All libraries are compiled using default compile parameters.

<sup>2</sup> The test inputs for libjpeg is taken from Google Image Test Suite.

<sup>3</sup> The number of test inputs in test suite is shown in the bracket.

Note this is to demonstrate the ability of inconsistent instructions on anti-fuzzing tasks. How to stealthily use these instructions is out of our scope. It is not easy for attackers to precisely recognize all the inconsistent instructions, which will be discussed in detail (Section 5).

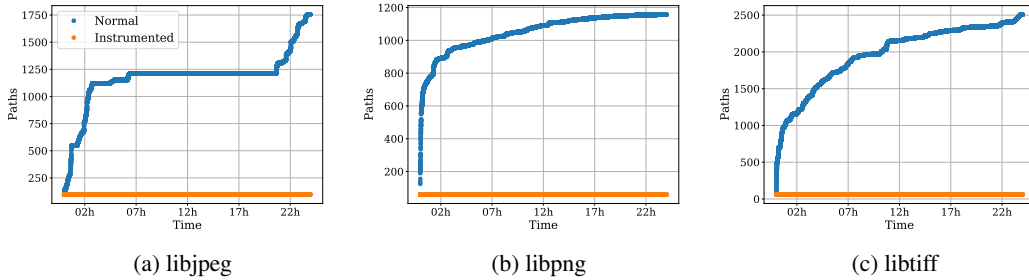


Figure 10: The result of Anti-Fuzzing experiment on three libraries. The blue lines show the coverage over 24 hours of fuzzing. The orange line shows the coverage for instrumented binaries, which decreases due to failed executions of QEMU.

**Answer to RQ3:** The inconsistent instructions are useful. We demonstrate that the inconsistent instructions can be used to detect the existence of the CPU emulator and prevent the malicious behavior from being monitored by dynamic analysis frameworks. Furthermore, the path coverage of programs fuzzed in emulators can be highly decreased with the help of inconsistent instructions.

## 5 Discussion

**Testing Instructions in Privileged Environments** Currently, the generated instruction streams are tested under unprivileged mode in both CPU emulators and real devices. Some instruction streams may have different execution results under privileged mode. For instance, instruction `WFI`, which results in a bug of QEMU user-mode, may not be an inconsistent instruction while executing in privileged mode. We plan to port INSDet to kernel-space for a more thorough testing.

**Testing Instruction Stream Sequences** INSDet now tests only one instruction stream each time during the differential testing. We can also test multiple instruction streams (instruction stream sequences) in the differential testing. The instruction stream sequences may trigger multiple system states and can test the decoding/executing logic of real devices and CPU emulators towards different state flags. How to design representative instruction stream sequences, and how to locate the inconsistent one will be the challenge. This will be our future work. Nevertheless, We have already discovered a huge number of inconsistent instruction streams with INSDet, covering 47.8% of the whole instructions. Every instruction stream sequence that contain the inconsistent instruction stream can result in inconsistent behaviors.

**Detecting (Ab)Used Inconsistent Instructions** The experiments in Section 4.3.2 and 4.3.3 show that attackers or vendors can (ab)use these inconsistent instructions to prevent the released binary from being monitored or fuzzed. It is not easy to recognize these inconsistent instructions due to the huge number of inconsistent instruction streams from 511 instructions. Some of these instructions are even commonly used (i.e., `STR` and `BLX` instruction). Apart from this, attackers can encrypt these instruction streams as data. Then these encrypted instruction streams can be decrypted and executed during runtime, which can increase the bar for detection. Thus, it is challenging to precisely detect those



inconsistent instruction streams if they are used to conduct anti-emulation or anti-fuzzing tasks. Furthermore, how to hide these inconsistent instruction streams from being detected is a *Cat and Mouse* problem. Stealthily using these instructions is not the purpose of our work.

**Other Emulators and Architectures** We test QEMU as it is one of the state-of-the-art emulators, which is well maintained and used by many industry tools (e.g., Android Studio) and academic prototypes [14, 18, 44, 9, 40, 10, 12, 21, 11, 16, 31]. However, INSDet can also be used to test the other emulators. Generally, all the tools (e.g., angr, Unicorn) including an emulated execution engine can be tested, which will be left as the future work. The whole framework of INSDet is architecture independent. However, we rely on ARM ASL to generate the test cases, which can explore multiple behaviors. If other architectures propose such kinds of specification language, we are able to generate the test cases. Otherwise, new test case generation algorithm should be developed.

## 6 Related Work

### 6.1 Testing CPU Emulators

Several works are proposed to test the CPU emulators. Lorenzo et al. proposed EmuFuzzer to test the CPU emulators [35, 34]. However, the seed used for testing mainly relies on randomization and a CPU-assisted mechanism, which may not cover all the CPU behaviors. Apart from testing user-level instructions, KEmuFuzzer is proposed to test the whole system emulators [33]. However, KEmuFuzzer relies on the manually written template to generate test cases. For better test case coverage, PokeEMU [32] is proposed. PokeEMU utilizes binary symbolic execution to generate more test cases from a high-fidelity emulator and apply these test cases on low-fidelity emulators. However, whether the high-fidelity emulator strictly follow the rule of specification is unknown. Furthermore, all the above mentioned works target on x86/x64 architectures. With the development of embedded systems and mobiles, the faithfully emulating ability for ARM architecture is a urgent need. Our work targets on ARM architecture and generates test cases from the specification itself (i.e., ARM ASL). The evaluation results show that we can find the real bugs and many inconsistent implementations between real devices and emulators, which can be abused by attackers.

### 6.2 Differential Testing

Differential testing is introduced by McKeeman et al. [36] to detect the implementation bugs by comparing the inconsistent behaviors between different software. For example, Yang et al. proposed Csmith, a powerful tool that can generate multiple C programs. With Csmith, hundreds of bugs are detected in the C compiler. Regarding the same goal, Le et al. introduced equivalence modulo inputs (EMI) [28] and many other differential testing tools are built based on EMI to validate the compiler implementations [29, 39].

Apart from testing compilers, researchers also utilize differential testing to validate the Database Management Systems (DBMS). Slutz et al. proposed the tool RAGS to explore bugs by executing different SQL queries on multiple DBMS. Though it is effective, it can only support a small set of SQL statements. Gu et al. evaluate the accuracy of DBMS optimizer by using options and hints to force the generation of different query plans. Jung et al. developed APOLLO [23] to test the performance regression bugs in DBMSs.

Furthermore, differential testing is powerful and applied to different domains such as testing SMT solvers [42, 41], JVM implementations [25], symbolic execution engines [25], and PDF readers [27].

### 6.3 Anti-Emulation Technique

Previous anti-emulation works [37] divide the anti-emulation technique into three categories. They are differences in behavior, differences in timing, and hardware specific values. Our work can automatically locate the inconsistent instructions, which result in different behavior and can be used by the previous anti-emulation technique. Jang et al. [20] address the importance of anti-emulation techniques on protecting the Commercial-Off-the-Shelf (COTS) software from being debugged or used without buying hardware. They propose three different anti-emulation techniques. However, some techniques rely on the race condition and are not easy to trigger.

## 7 Conclusion

We design and implement INSDet, a framework that can automatically locate the inconsistent ARM instructions, which can result in inconsistent behaviors between CPU emulator and real devices. With INSDet, we generate 2,774,649 representative instruction streams and detect 155,642 inconsistent instruction streams covering 30% of the

instruction encodings and 47.8% instructions. By analyzing the root cause of inconsistent instructions, we noticed four bugs of QEMU, which are confirmed by QEMU developers, covering 13 instruction encodings including very commonly used ones (e.g., STR, BLX). We also demonstrate the capability of inconsistent instructions on detecting emulators, anti-emulation, and anti-fuzzing.

## References

- [1] 64 bit junos arm® development platform. <https://developer.arm.com/-/media/Arm%20Developer%20Community/PDF/Juno%20r2%20datasheet.pdf>.
- [2] Afl qemu mode: high-performance binary-only instrumentation for afl-fuzz. [https://github.com/google/afl/tree/master/qemu\\_mode](https://github.com/google/afl/tree/master/qemu_mode).
- [3] ARM Exploration tools. <https://developer.arm.com/architectures/cpu-architecture/a-profile/exploration-tools>.
- [4] Capstone. <https://www.capstone-engine.org/>.
- [5] Panda.re. <https://panda.re/>.
- [6] Suterusu. <https://github.com/mncoppola/suterusu>.
- [7] TriforceAFL. <https://github.com/nccgroup/TriforceAFL>.
- [8] Z3Prover. <https://github.com/Z3Prover/z3>.
- [9] Abdulla Alwabel, Hao Shi, Genevieve Bartlett, and Jelena Mirkovic. Safe and automated live malware experimentation on public testbeds. In *Proceedings of the 7th Workshop on Cyber Security Experimentation and Test*, 2014.
- [10] Curtis Carmony, Xunchao Hu, Heng Yin, Abhishek Vasisht Bhaskar, and Mu Zhang. Extract me if you can: Abusing pdf parsers in malware detectors. In *Proceedings of the 23rd Annual Network and Distributed System Security Symposium*, 2016.
- [11] Daming D Chen, Maverick Woo, David Brumley, and Manuel Egele. Towards automated dynamic analysis for linux-based embedded firmware. In *Proceedings of the 23rd Annual Network and Distributed System Security Symposium*, 2016.
- [12] Vitaly Chipounov, Volodymyr Kuznetsov, and George Candea. S2e: A platform for in-vivo multi-path analysis of software systems. *Acm Sigplan Notices*, 2011.
- [13] Abraham A Clements, Eric Gustafson, Tobias Scharnowski, Paul Grosen, David Fritz, Christopher Kruegel, Giovanni Vigna, Saurabh Bagchi, and Mathias Payer. Halucinator: Firmware re-hosting through abstraction layer emulation. In *Proceedings of the 29th USENIX Security Symposium*, 2020.
- [14] Ali Davanian, Zhenxiao Qi, Yu Qu, and Heng Yin. Decaf++: Elastic whole-system dynamic taint analysis. In *Proceedings of the 22nd International Symposium on Research in Attacks, Intrusions and Defenses*, 2019.
- [15] Bo Feng, Alejandro Mera, and Long Lu. P2IM: Scalable and hardware-independent firmware testing via automatic peripheral interface modeling. In *Proceedings of the 29th USENIX Security Symposium*, 2019.
- [16] Qian Feng, Aravind Prakash, Heng Yin, and Zhiqiang Lin. Mace: High-coverage and robust memory analysis for commodity operating systems. In *Proceedings of the 30th annual computer security applications conference*, 2014.
- [17] Emre Güler, Cornelius Aschermann, Ali Abbasi, and Thorsten Holz. Antifuzz: Impeding fuzzing audits of binary executables. In *Proceedings of the 28th {USENIX} Security Symposium*, 2019.
- [18] Andrew Henderson, Aravind Prakash, Lok Kwong Yan, Xunchao Hu, Xujiwen Wang, Rundong Zhou, and Heng Yin. Make it work, make it right, make it fast: building a platform-neutral whole-system dynamic binary analysis platform. In *Proceedings of the 2014 International Symposium on Software Testing and Analysis*, 2014.
- [19] Anoirel Issa. Anti-virtual machines and emulations. *Journal in Computer Virology*, 2012.
- [20] Daehee Jang, Yunjong Jeong, Sungman Lee, Minjoon Park, Kuenhwan Kwak, Donguk Kim, and Brent Byunghoon Kang. Rethinking anti-emulation techniques for large-scale software deployment. *Computers & Security*, 2019.
- [21] Xuxian Jiang, Xinyuan Wang, and Dongyan Xu. Stealthy malware detection and monitoring through vmm-based “out-of-the-box” semantic view reconstruction. *ACM Transactions on Information and System Security*, 2010.

- [22] Evan Johnson, Maxwell Bland, YiFei Zhu, Joshua Mason, Stephen Checkoway, Stefan Savage, and Kirill Levchenko. Jetset: Targeted firmware rehosting for embedded systems. In *Proceedings of the 30th {USENIX} Security Symposium*, 2021.
- [23] Jinho Jung, Hong Hu, Joy Arulraj, Taesoo Kim, and Woonhak Kang. Apollo: Automatic detection and diagnosis of performance regressions in database systems. *Proceedings of the VLDB Endowment*, 2019.
- [24] Jinho Jung, Hong Hu, David Solodukhin, Daniel Pagan, Kyu Hyung Lee, and Taesoo Kim. Fuzzification: Anti-fuzzing techniques. In *Proceedings of the 28th {USENIX} Security Symposium*, 2019.
- [25] Timotej Kapus and Cristian Cadar. Automatic testing of symbolic execution engines via program generation and differential testing. In *Proceedings of the 32nd IEEE/ACM International Conference on Automated Software Engineering*, 2017.
- [26] Mingeun Kim, Dongkwan Kim, Eunsoo Kim, Suryeon Kim, Yeongjin Jang, and Yongdae Kim. Firmae: Towards large-scale emulation of iot firmware for dynamic analysis. In *Proceedings of the 2020 Annual Computer Security Applications Conference*, 2020.
- [27] Tomasz Kuchta, Thibaud Lutellier, Edmund Wong, Lin Tan, and Cristian Cadar. On the correctness of electronic documents: studying, finding, and localizing inconsistency bugs in pdf readers and files. *Empirical Software Engineering*, 2018.
- [28] Vu Le, Mehrdad Afshari, and Zhendong Su. Compiler validation via equivalence modulo inputs. *ACM SIGPLAN Notices*, 2014.
- [29] Christopher Lidbury, Andrei Lascu, Nathan Chong, and Alastair F Donaldson. Many-core compiler fuzzing. *ACM SIGPLAN Notices*, 2015.
- [30] Cătălin Valeriu Liță, Doina Cosovan, and Dragoș Gavriluț. Anti-emulation trends in modern packers: a survey on the evolution of anti-emulation techniques in upa packers. *Journal of Computer Virology and Hacking Techniques*, 2018.
- [31] Lannan Luo, Yu Fu, Dinghao Wu, Sencun Zhu, and Peng Liu. Repackage-proofing android apps. In *Proceedings of the 46th Annual IEEE/IFIP International Conference on Dependable Systems and Networks*, 2016.
- [32] Lorenzo Martignoni, Stephen McCamant, Pongsin Poosankam, Dawn Song, and Petros Maniatis. Path-exploration lifting: Hi-fi tests for lo-fi emulators. In *Proceedings of the Seventeenth International Conference on Architectural Support for Programming Languages and Operating Systems*, 2012.
- [33] Lorenzo Martignoni, Roberto Paleari, Giampaolo Fresi Roglia, and Danilo Bruschi. Testing system virtual machines. In *Proceedings of the 19th international symposium on software testing and analysis*, 2010.
- [34] Lorenzo Martignoni, Roberto Paleari, Alessandro Reina, Giampaolo Fresi Roglia, and Danilo Bruschi. A methodology for testing cpu emulators. *ACM Transactions on Software Engineering and Methodology*, 2013.
- [35] Lorenzo Martignoni, Roberto Paleari, Giampaolo Fresi Roglia, and Danilo Bruschi. Testing cpu emulators. In *Proceedings of the eighteenth international symposium on Software testing and analysis*, 2009.
- [36] William M McKeeman. Differential testing for software. *Digital Technical Journal*, 1998.
- [37] Thomas Raffetseder, Christopher Kruegel, and Engin Kirda. Detecting system emulators. In *Proceedings of the 2007 International Conference on Information Security*. Springer, 2007.
- [38] Alastair Reid. Trustworthy specifications of arm@ v8-a and v8-m system level architecture. In *Proceedings of the 16th Formal Methods in Computer-Aided Design*, 2016.
- [39] Chengnian Sun, Vu Le, and Zhendong Su. Finding compiler bugs via live code mutation. In *Proceedings of the 2016 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications*, 2016.
- [40] Jinpeng Wei, Lok K Yan, and Muhammad Azizul Hakim. Mose: Live migration based on-the-fly software emulation. In *Proceedings of the 31st Annual Computer Security Applications Conference*, 2015.
- [41] Dominik Winterer, Chengyu Zhang, and Zhendong Su. On the unusual effectiveness of type-aware operator mutations for testing smt solvers. *Proceedings of the ACM on Programming Languages*, (OOPSLA), 2020.
- [42] Dominik Winterer, Chengyu Zhang, and Zhendong Su. Validating smt solvers via semantic fusion. In *Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation*, 2020.
- [43] Xu Chen, J. Andersen, Z. M. Mao, M. Bailey, and J. Nazario. Towards an understanding of anti-virtualization and anti-debugging behavior in modern malware. In *Proceedings of the 38th IEEE International Conference on Dependable Systems and Networks*, 2008.

- [44] Lok Kwong Yan and Heng Yin. Droidscape: Seamlessly reconstructing the {OS} and dalvik semantic views for dynamic android malware analysis. In *Proceedings of the 21st {USENIX} Security Symposium*, 2012.
- [45] Yaowen Zheng, Ali Davanian, Heng Yin, Chengyu Song, Hongsong Zhu, and Limin Sun. Firm-af: high-throughput greybox fuzzing of iot firmware via augmented process emulation. In *Proceedings of the 28th {USENIX} Security Symposium*, 2019.