# Symbolic Execution of x86 assembly in Isabelle/HOL

## Freek Verbeek
Virginia Tech, USA

freek@vt.edu

## Abhijith Bharadwaj
Virginia Tech, USA

## Joshua A. Bockenek
Virginia Tech, USA

## Ian Roessle
Joint Artificial Intelligence Center, Washington DC, US

## Binoy Ravindran
Virginia Tech, USA

──── **Abstract** ────

In this short paper we present progress on a symbolic execution engine for x86 assembly in the Isabelle/HOL theorem prover. We discuss the two main challenges tackled: 1.) how to leverage reliable machine-learned semantics of x86 assembly instructions, and 2.) how to generate preconditions that allow deterministic symbolic execution of basic blocks. We end with a discussion on how we intend to use our symbolic execution engine.

## 1 Introduction

Symbolic execution is a powerful technique in program verification and analysis [7, 2]. It can be used to explore an overapproximation of all possible paths. In case of assembly code, it can also be used to *summarize* state changes induced by sequences of individual assembly instructions. In assembly, one will typically find series of instructions whose net effect can be described much more succinctly than by using the semantics of the individual instructions. As example, consider the x86 assembly sequence `push rbp; pop rbp`. The net effect is only a single write into memory (register `rbp` is written to the top of the stack frame). The succinct output of symbolic execution can be the base for further for formal verification.

This short paper describes our progress in building a formal symbolic execution engine in Isabelle/HOL [5] for x86-64. Our symbolic execution engine targets basic blocks, i.e., blocks without unconditional jumps. We specifically deal with the following two challenges:

- The semantics of x86 are typically highly complicated and its CISC nature requires formal semantics for many instructions. The Intel manuals provide documentation, but translating these into a formal model is error-prone and requires human interpretation. We use *Strata* [4] to embed highly trustworthy machine-learned instruction semantics into Isabelle/HOL. The challenge is that, since these semantics are not manually written but machine-learned, they are typically not in a form suitable for formal verification. We

43 thus provide manually written *presimplified semantics* and prove equivalence between our
44 manually written semantics and the machine-learned version.

45 ▪ Once each instruction in a block has been given semantics, symbolic execution amounts
46 to aggregating these individual state changes. The objective is that one basic block has a
47 deterministic aggregated state change, since each individual instruction is deterministic.
48 However, since all values are symbolic, addresses are typically symbolic as well. This
49 leads to the *memory aliasing* problem: if two values are written to memory to symbolic
50 addresses $a_0$ and $a_1$, it is possible that they overwrite each other, overlap each other,
51 or are separate. We generate preconditions under which symbolic execution becomes
52 deterministic for basic blocks.

## 2 Using machine-learned semantics

54 Strata uses a stochastic search methodology to derive instruction semantics from an x86-
55 64 machine. The search space used to learn instructions consists of 62 hard-coded base
56 instructions. These base instructions cover bit-vector operations such as integer arithmetic,
57 bitwise operations, data movement, floating point operations, splitting and combining of
58 registers, and setting and clearing of status flags. The base set covers fundamental operations,
59 serving as building blocks for the more complex instructions. Ultimately, semantics are
60 learned as assignments of bit-vector formula's to state parts.

61 In [6], we describe a methodology for generalizing the output of Strata, and lifting it into
62 the Isabelle/HOL theorem prover. As an example, we consider the instruction variant `sub`
63 `r32 m32`, which subtracts the value stored in the 32-bit memory location from the value
64 stored in the 32-bit register. Note that in x86-64, a 32-bit register is actually the lower part
65 of a 64-bit register. This instruction thus actually reads from and writes to a 64-bit register.
66 We also show two of the flags: the zero flag and the carry flag.

$$
\begin{aligned}
r64 &\coloneqq \underset{32}{0} \smile \langle 31,0\rangle (\underset{1}{0} \smile \neg m32 + \underset{33}{1} + \underset{1}{0} \smile \langle 31,0\rangle (r64)) \\
ZF &\coloneqq \langle 31,0\rangle (\underset{1}{0} \smile \neg m32 + \underset{33}{1} + \underset{1}{0} \smile \langle 31,0\rangle (r64)) == \underset{32}{0} \\
CF &\coloneqq \langle 32,32\rangle (\underset{1}{0} \smile \neg m32 + \underset{33}{1} + \underset{1}{0} \smile \langle 31,0\rangle (r64)) == \underset{32}{1}
\end{aligned}
$$

68 It can be seen that the semantics are expressed in base instructions such as concatenation
69 ($\smile$), taking a sub-bit-vector ($\langle 31,0\rangle$), negation, addition, constants ($\underset{33}{1}$ means "the constant
70 1 in 33-bit mode) and equality. These semantics, however, also seem overly complicated. In
71 order to express the semantics of the zero flag, for example, the input values are extended to
72 33-bit mode, after which a two's complement subtraction happens. Then the lower 32 bits of
73 33 are compared to 0. Humanly defined semantics would simply state `r32 == m32`, i.e., the
74 zero flag after subtraction is set when its inputs are equal. We thus defined the following
75 manually written presimplified semantics:

$$
\begin{aligned}
r64 &\coloneqq \text{zextend}(\langle 31,0\rangle(r64) - m32) \\
ZF &\coloneqq \langle 31,0\rangle r64 = m32 \\
CF &\coloneqq \langle 31,0\rangle r64 < m32
\end{aligned}
$$

79 These two semantics are formally proven to be equivalent. We have presimplified semantics
80 for 84 instruction mnemonics, where each mnemonic has several variants. For example, `sub`
81 is a mnemonic with 8, 16, 32 and 64 bit variants, and each of these variants has further
82 variation in whether its operands are memory or registers. Not all semantics come from
83 Strata, e.g., the shift instructions have been defined manually. More details can be found
84 in [6].

## 3    Determinizing Symbolic Execution

Consider the following x86 assembly sequence:

```
mov QWORD PTR [rsp-16], 1
mov DWORD PTR [rsp-24], 2
mov rax, QWORD PTR [rsp-16]
```

The first instruction moves the quad (8 byte) word `1` to memory location `[rsp-16]`. The second moves the 4 byte word `2` to memory location `[rsp-24]`. The third moves 8 bytes from memory location `[rsp-16]` into the `RAX` register. Symbolic execution should produce the following:

$$s' = s([\mathtt{rsp} - 16] := \underset{64}{1}, [\mathtt{rsp} - 24] := \underset{32}{2}, \mathtt{RAX} := \underset{64}{1})$$

That is, the new state $s'$ is the result of three state changes with respect to the input state $s$. For sake of presentation, the instruction pointer is omitted. We illustrate that in order to get this seemingly trivial result, we require both extra preconditions and solving of linear equations.

Symbolic execution starts in state $s$ and sequentially applies the presimplified semantics of the current instruction. After execution of the first instruction, the current symbolic state is:

$$s' = s([\mathtt{rsp} - 16] := \underset{64}{1})$$

Now, in order to execute the second instruction, it needs to be established that the two regions written to are separate. If they are separate, the next symbolic state is equal to:

$$s' = s([\mathtt{rsp} - 16] := \underset{64}{1}, [\mathtt{rsp} - 24] := \underset{32}{2})$$

However, were they to overlap, then a different symbolic state would be produced.

To know whether they are separate, the following linear equation must be solved:

$$\mathtt{rsp} - 16 + 8 \le \mathtt{rsp} - 24 \lor \mathtt{rsp} - 24 + 4 \le \mathtt{rsp} - 16$$

This seems a trivial linear equation, since $\mathtt{rsp} - 20 \le \mathtt{rsp} - 16$. However, the addresses are computed in 64-bit mode, i.e., the address computations are modulo $2^{64}$. Thus, the equation is not true: if for example $\mathtt{rsp} = 16$, then $\mathtt{rsp} - 20 > \mathtt{rsp} - 16$. When the extra precondition $\mathtt{rsp} \ge 20$ is assumed, the linear equation can be solved and we can complete symbolic execution deterministically.

Our solution to this problem is as follows:

1. For each basic block, identify the accessed regions;
2. For each region, generate the preconditions necessary to prevent under- and overflow;
3. For each pair of two regions, precompute whether they are separate, and whether they are enclosed in each other;
4. For each basic block, generate a lemma in Isabelle/HOL with as assumptions the generated preconditions and the precomputed relations.

Step 3 uses the Z3 theorem prover [3]: for each pair of regions, linear equations are generated that model separation and enclosure. This also prevents a vacuous truth: since the assumptions are generated, we need to make sure that they are internally consistent. Z3 ensures that we cannot add an assumption such as "$[\mathtt{rsp} - 16, 8]$ is separate from $[\mathtt{rsp} - 12, 8]$".

Note that in the given example, the basic block *is* deterministic. In general, that is not necessarily the case, e.g., in case of aliasing. Consider the following example:

```
127  mov QWORD PTR [rdi], 1
128  mov DWORD PTR [rsi], 2
129  mov rax, QWORD PTR [rdi]
```

Symbolic execution cannot produce a deterministic value for register `RAX`, since it depends on the values of registers `RDI` and `RSI`. In this case, the Isabelle symbolic execution will fail, and the user manually needs to insert as assumption that, e.g., $[RDI, 8]$ is separate from $[RSI, 4]$. More details can be found in [1].

## 4    Use Cases of Formal Symbolic Execution

As conclusion, we discuss some use cases of formal symbolic execution.

**Combine with CFG extraction** As discussed, we do symbolic execution per basic block. The CFG dictates how these basic blocks are tied together. We aim to combine a formally proven correct CFG extraction tool with our symbolic execution engine to get a summarized – but correct – representation of the binary in the theorem prover.

**Formal Proofs of Memory Usage** In [1], we use the symbolic execution engine to reason over the memory read from and written to by functions in binaries. We generate Floyd invariants, that allow reasoning per basic block to be used to reason over a function as a whole. We have applied this to 71 functions of the binary of HermitCore, and dealt with functions with loops, recursion, and pointer arguments. The methodology requires interactive theorem proving, and we aim to make this methodology more automatable to achieve better scalability.

**Formal Proofs of Soundness of Randomizers** Binary randomization is a technique used to prevent return-oriented-programming attacks. A randomizer rewrites basic blocks to eliminate so-called gadgets, i.e., byte-sequences that can be interpreted as a `ret` instruction. Analysis of these randomizers typically focuses on security properties, and less on soundness. Using formal symbolic execution, we intend to compare the semantics of a binary with the semantics of its randomized version, and thereby prove soundness.

### ── References ──

**1**    Joshua A. Bockenek, Freek Verbeek, Peter Lammich, and Binoy Ravindran. Formal verification of memory preservation of x86-64 binaries. In *International Conference on Computer Safety, Reliability and Security (SAFECOMP'19)*, 2019. To be published.

**2**    Eric Cheng. *Binary Analysis and Symbolic Execution with angr*. PhD thesis, The MITRE Corporation, 2016.

**3**    Leonardo De Moura and Nikolaj Bjørner. Z3: An efficient smt solver. In *International conference on Tools and Algorithms for the Construction and Analysis of Systems*, pages 337–340. Springer, 2008.

**4**    Stefan Heule, Eric Schkufza, Rahul Sharma, and Alex Aiken. Stratified synthesis: automatically learning the x86-64 instruction set. In *ACM SIGPLAN Notices*, volume 51, pages 237–250. ACM, 2016.

**5**    Tobias Nipkow, Lawrence C Paulson, and Markus Wenzel. *Isabelle/HOL: a proof assistant for higher-order logic*, volume 2283. Springer Science & Business Media, 2002.

**6**    Ian Roessle, Freek Verbeek, and Binoy Ravindran. Formally verified big step semantics out of x86-64 binaries. In *Proceedings of CPP'19*, pages 181–195. ACM, 2019.

**7**    Teodor Stoenescu, Alin Stefanescu, Sorina Predut, and Florentin Ipate. River: A binary analysis framework using symbolic execution and reversible x86 instructions. In *International Symposium on Formal Methods*, pages 779–785. Springer, 2016.