

Verified Abstract Interpretation Techniques for Disassembling Low-level Self-modifying Code

Sandrine Blazy¹ Vincent Laporte¹ David Pichardie²

¹Université Rennes 1 – IRISA – Inria

²ENS Rennes – IRISA – Inria

July 16th 2014

Static analysis. . .

Prove program safety before running it

. . . of binary

Source not available

Compiler not trusted

Self-Modifying

Packed software

Obfuscation

JIT compiler

```
07 00 06 07 03 00 00 00 00 00 00 05
00 00 00 00 00 00 01 00 09 00 00 00
00 00 00 04 09 00 00 02 00 00 00 02
05 00 00 02 04 00 00 00 00 00 00 01
?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ??
?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ??
```

```
R0: ?? ?? ?? ??
R1: ?? ?? ?? ??
R2: ?? ?? ?? ??
R3: ?? ?? ?? ??
R4: ?? ?? ?? ??
R5: ?? ?? ?? ??
R6: 04 a4 3b 09
R7: 10 b9 e6 51
Flags:      EQ
```

Next instruction: `cmp R6, R7`

```
07 00 06 07 03 00 00 00 00 00 00 05
00 00 00 00 00 00 01 00 09 00 00 00
00 00 00 04 09 00 00 02 00 00 00 02
05 00 00 02 04 00 00 00 00 00 00 01
?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ??
?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ??
```

```
R0: ?? ?? ?? ??
R1: ?? ?? ?? ??
R2: ?? ?? ?? ??
R3: ?? ?? ?? ??
R4: ?? ?? ?? ??
R5: ?? ?? ?? ??
R6: 04 a4 3b 09
R7: 10 b9 e6 51
Flags: LT LE
```

Next instruction: gotoLE 5

```
07 00 06 07 03 00 00 00 00 00 00 05
00 00 00 00 00 00 01 00 09 00 00 00
00 00 00 04 09 00 00 02 00 00 00 02
05 00 00 02 04 00 00 00 00 00 00 01
?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ??
?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ??
```

```
R0: ?? ?? ?? ??
R1: ?? ?? ?? ??
R2: ?? ?? ?? ??
R3: ?? ?? ?? ??
R4: ?? ?? ?? ??
R5: ?? ?? ?? ??
R6: 04 a4 3b 09
R7: 10 b9 e6 51
Flags: LT LE
```

Next instruction: `cst 4` → R₀

```
07 00 06 07 03 00 00 00 00 00 00 05
00 00 00 00 00 00 01 00 09 00 00 00
00 00 00 04 09 00 00 02 00 00 00 02
05 00 00 02 04 00 00 00 00 00 00 01
?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ??
?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ??
```

```
R0: 00 00 00 04
R1: ?? ?? ?? ??
R2: ?? ?? ?? ??
R3: ?? ?? ?? ??
R4: ?? ?? ?? ??
R5: ?? ?? ?? ??
R6: 04 a4 3b 09
R7: 10 b9 e6 51
Flags: LT LE
```

Next instruction: `cst 2` → R₂

```
07 00 06 07 03 00 00 00 00 00 00 05
00 00 00 00 00 00 01 00 09 00 00 00
00 00 00 04 09 00 00 02 00 00 00 02
05 00 00 02 04 00 00 00 00 00 00 01
?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ??
?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ??
```

```
R0: 00 00 00 04
R1: ?? ?? ?? ??
R2: 00 00 00 02
R3: ?? ?? ?? ??
R4: ?? ?? ?? ??
R5: ?? ?? ?? ??
R6: 04 a4 3b 09
R7: 10 b9 e6 51
Flags: LT LE
```

Next instruction: store R₀ → *R₂

```
07 00 06 07 03 00 00 00 00 00 00 04
00 00 00 00 00 00 01 00 09 00 00 00
00 00 00 04 09 00 00 02 00 00 00 02
05 00 00 02 04 00 00 00 00 00 00 01
?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ??
?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ??
```

```
R0: 00 00 00 04
R1: ?? ?? ?? ??
R2: 00 00 00 02
R3: ?? ?? ?? ??
R4: ?? ?? ?? ??
R5: ?? ?? ?? ??
R6: 04 a4 3b 09
R7: 10 b9 e6 51
Flags: LT LE
```

Next instruction: goto 1


```
07 00 06 07 03 00 00 00 00 00 00 04
00 00 00 00 00 00 01 00 09 00 00 00
00 00 00 04 09 00 00 02 00 00 00 02
05 00 00 02 04 00 00 00 00 00 00 01
?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ??
?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ??
```

```
R0: 00 00 00 04
R1: ?? ?? ?? ??
R2: 00 00 00 02
R3: ?? ?? ?? ??
R4: ?? ?? ?? ??
R5: ?? ?? ?? ??
R6: 04 a4 3b 09
R7: 10 b9 e6 51
Flags: LT LE
```

Next instruction: gotoLE 4

```
07 00 06 07 03 00 00 00 00 00 00 04
00 00 00 00 00 00 01 00 09 00 00 00
00 00 00 04 09 00 00 02 00 00 00 02
05 00 00 02 04 00 00 00 00 00 00 01
?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ??
?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ??
```

```
R0: 00 00 00 04
R1: ?? ?? ?? ??
R2: 00 00 00 02
R3: ?? ?? ?? ??
R4: ?? ?? ?? ??
R5: ?? ?? ?? ??
R6: 04 a4 3b 09
R7: 10 b9 e6 51
Flags: LT LE
```

Next instruction: `halt R0`

```
07 00 06 07 03 00 00 00 00 00 00 04
00 00 00 00 00 00 01 00 09 00 00 00
00 00 00 04 09 00 00 02 00 00 00 02
05 00 00 02 04 00 00 00 00 00 00 01
?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ??
?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ?? ??
```

```
R0: 00 00 00 04
R1: ?? ?? ?? ??
R2: 00 00 00 02
R3: ?? ?? ?? ??
R4: ?? ?? ?? ??
R5: ?? ?? ?? ??
R6: 04 a4 3b 09
R7: 10 b9 e6 51
Flags: LT LE
```

Final value: 4

Given such a program, how to...

Disassemble it?

Where are the instructions?

What are they?

Compute its control-flow graph?

What are the targets of the computed jumps?

Automatically prove safety properties about it?

Can its execution be stuck?

May it access secret parts of the memory?

Trust the answers to above questions?

Is the analysis sound?

“Certified Self-Modifying Code”

(Cai, Shao, and Vaynberg PLDI'2007)

- Framework for manual verification of self-modifying programs

“WYSINWYX: What You See Is Not What You eXecute”

(Balakrishnan and Reps 2010)

- Static analysis of x86

Formalize in Coq a static analysis

- that is flow sensitive
 - attach to each reachable program point an over-approximation of the state at that point
 - analyze the content of the memory and of the registers
- without a previous disassembling or CFG reconstruction

Abstract interpreter

State abstraction

Numeric abstraction

Each layer is parameterized by the underlying one

Numeric abstract domains

Abstract sets of 32-bit machine integers

Finite sets e.g., $\{0; 1; 7\}$

Strided Intervals that combines interval and congruence information
e.g., $[1000; 2000].4$ represents $\{1000; 10004; 1008; \dots; 2000\}$

Signature of abstract domains

Each abstract domain is equipped with a lattice structure

```
Class weak_lattice (A: Type) : Type :=  
{ leb: A → A → bool      (* order *)  
; top: A                    (* maximal element *)  
; join: A → A → A         (* binary upper bound *)  
; widen: A → A → A        (* extrapolation operator *)  
}.
```

Concretization based specification

Each abstract domain comes with a **concretization relation**

```
Class gamma_op (A B: Type) : Type :=  $\gamma : A \rightarrow \mathcal{P}(B)$ .
```

```
Record adom (A B:Type)(_:weak_lattice A)(_:gamma_op A B): Prop :=  
{ gamma_monotone:  $\forall a1 a2, \text{leb } a1 a2 = \text{true} \rightarrow \gamma a1 \subseteq \gamma a2$   
; gamma_top:  $\forall x, x \in \gamma \text{ top}$   
; join_sound:  $\forall x y, \gamma x \cup \gamma y \subseteq \gamma (\text{join } x y)$  }.
```

- No more properties required

Example (Strided intervals)

```
Instance si_gamma : gamma_op strided_interval int :=  $\lambda x i,$   
  low_bound x  $\leq$  Int.signed i  $\leq$  up_bound x  
   $\wedge$  low_bound x  $\equiv$  Int.signed i [ stride x ].
```

Signature of the state abstraction

```
Record mem_dom (ab_mem: Type) (ab_num: Type) : Type :=  
{ (* lattice *)  
  as_wl: weak_lattice ab_mem  
  (* queries *)  
; var: ab_mem → register → ab_num  
; load: ab_mem → addr → ab_num  
  (* abstract transformers *)  
; assign: ab_mem → register → ab_num → ab_mem  
  (* more omitted ... *)  
}.
```

Signature of the state abstraction

```
Record mem_dom (ab_mem: Type) (ab_num: Type) : Type :=
{ (* lattice *)
  as_wl: weak_lattice ab_mem
  (* queries *)
; var: ab_mem → register → ab_num
; load: ab_mem → addr → ab_num
  (* abstract transformers *)
; assign: ab_mem → register → ab_num → ab_mem
  (* more omitted ... *)
}.
```

Parameterized by a numeric abstraction: `ab_num`

Signature of the state abstraction

```
Record mem_dom (ab_mem: Type) (ab_num: Type) : Type :=  
{ (* lattice *)  
  as_wl: weak_lattice ab_mem  
  (* queries *)  
; var: ab_mem → register → ab_num  
; load: ab_mem → addr → ab_num  
  (* abstract transformers *)  
; assign: ab_mem → register → ab_num → ab_mem  
  (* more omitted ... *)  
}.
```

Equipped with a lattice structure: order, top element, join operator

Signature of the state abstraction

```
Record mem_dom (ab_mem: Type) (ab_num: Type) : Type :=
{ (* lattice *)
  as_wl: weak_lattice ab_mem
  (* queries *)
; var: ab_mem → register → ab_num
; load: ab_mem → addr → ab_num
  (* abstract transformers *)
; assign: ab_mem → register → ab_num → ab_mem
  (* more omitted ... *)
}.
```

Provides access to the abstraction of each memory unit

Signature of the state abstraction

```
Record mem_dom (ab_mem: Type) (ab_num: Type) : Type :=
{ (* lattice *)
  as_wl: weak_lattice ab_mem
  (* queries *)
; var: ab_mem → register → ab_num
; load: ab_mem → addr → ab_num
  (* abstract transformers *)
; assign: ab_mem → register → ab_num → ab_mem
  (* more omitted ... *)
}.
```

Abstract operators model concrete instructions

Signature of the state abstraction

```
Record mem_dom (ab_mem: Type) (ab_num: Type) : Type :=  
{ (* lattice *)  
  as_wl: weak_lattice ab_mem  
  (* queries *)  
; var: ab_mem → register → ab_num  
; load: ab_mem → addr → ab_num  
  (* abstract transformers *)  
; assign: ab_mem → register → ab_num → ab_mem  
  (* more omitted ... *)  
}.
```

Example (Specification of the load query)

$$\text{load_sound: } \forall \text{ ab, } \forall \text{ m,}$$
$$\text{m} \in \gamma(\text{ab}) \rightarrow \forall \text{ a:addr, m(a) } \in \gamma(\text{load ab a})$$

Abstract small-step

Specification

Lemma `ab_step_correct` : $\forall m m' ab,$
 $m \in \gamma(ab) \rightarrow m \rightsquigarrow m' \rightarrow m' \in \gamma(ab_step\ m.\ (pp)\ ab)$

Algorithm

Given an abstract state `m` at an address `pp` (program point)

- 1 Decode all possible instructions from `pp`
- 2 For each of them predict
 - next program point
 - next abstract state
- 3 Propagate the results

Fix-point computation

Given a program P (partial initial memory), compute an abstract environment E such that

- $\text{init}(P) \sqsubseteq E[0]$, and
- $\forall p \ p' \ m', (p', m') \in \text{ab_step}(p)(E[p]) \rightarrow m' \sqsubseteq E[p']$

- Iterate from initial state
- Work-set of reachable program points that need further analysis
- Widening steps to ensure termination
 - The widening policy is a parameter of the analysis
- A posteriori validation

Lemma `validate_correct` : $\forall P \ E,$
`validate_fixpoint` $P \ E = \text{true} \rightarrow \llbracket P \rrbracket \sqsubseteq \gamma(E)$

What can we do with this fix-point?

Given a (sound) fix-point:

Program safety can be proved (no run-time error)

Theorem `analysis_sound` : $\forall P \text{ dom fuel ab_num},$
`analysis ab_num P dom fuel` $\neq \text{None} \rightarrow \text{safe } P$

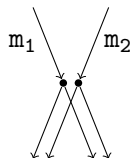
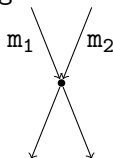
A CFG can be computed

- Reachable program points
- Possible instructions when execution reaches a program point
- Possible targets of the jumps

Extension: Trace Partitioning

Some merges of abstract states incur a large precision loss:

- Do not merge abstract states that differ according to some criterion



- Generalizes the flow sensitivity
- Need to modify only the fix-point computation (and checker)

Experimental evaluation

Extract to OCaml and run on challenging examples¹

- ✓ **opcode modification** overwrites code before execution
- ✓ **multilevel RCG** code that once executed produces code that once. . .
- × **bootloader** loads code from disk
- ✓ **control flow modification** overwrites jumps
- ✓ **vector dot product** specializes a program for a given argument
- ✓ **runtime code checking** validates code integrity
- ✓ **Fibonacci sequence** uses an instruction as accumulator
- × **self-replication** infinitely copies itself
- ✓ **mutual modification** two parts modify each other
- ✓ **polymorphic code** mutates itself to escape anti-viruses
- ✓ **code obfuscation** features fake instructions if naively disassembled
- ✓ **code encryption** decrypts a code, runs it, and crypts it back

¹Cai, Shao, and Vaynberg PLDI'2007.

Conclusion

Summary

- Verified executable abstract interpreter for a binary language
- Handles seamlessly self-modifying programs
- Able to prove the safety of challenging examples
- Full development available online:
<http://www.irisa.fr/celtique/ext/smc/>

Perspectives

- More realistic languages (x86)
- Generate the analyzer (and its proof) to cope with the huge number of instructions