### Université BORDEAUX

## **Binary Code Analysis: Concepts and Perspectives**

# Emmanuel Fleury <emmanuel.fleury@u-bordeaux.fr>

LaBRI, Université de Bordeaux, France

May 12, 2016



**1** Introducing to Binary Code Analysis

- 2 Why Is Binary Analysis Special?
- **3** Low-level Programs Formal Model
- 4 Control-flow Recovery
- 5 Current and Future Trends



#### Introducing to Binary Code Analysis

- Basic Definitions
- Binary Analysis Pipeline
- Practical and Theoretical Challenges

#### 2 Why Is Binary Analysis Special?

- 3 Low-level Programs Formal Model
- 4 Control-flow Recovery
- **5** Current and Future Trends

- Analysis of legacy/off-the-shelf/proprietary software;
- Software reverse-engineering on malware (or others);
- Analysis of software generated with untrusted compiler;
- To capture many low-level security issues;
- Analysis of low-level interactions (hardware/OS).
- Optimize a binary without the sources (recompilation).

## What we mean by "Binary Programs"? Université BORDEAUX

**Abstract Model**: All **unnecessary information** for the analysis have been removed. Only **necessary information** remains.

**Source Code**: Keep track of high-level information about the program such as variables, types, functions. But also, variable and function names, and pragmas or code decorations.

**Bytecode**: May vary depending on the bytecode considered, but keep track of few high-level information about the program such as **types** and **functions**. But, programs are usually **unstructured**.

**Binary File**: Only keep track of the **instructions** in an **unstructured way** (no forloop, no clear argument passing in procedures, ...). No type, no naming. But, the binary file may enclose **meta-data** that might be helpful (symbols, debug, ...).

**Memory Dump**: Pure assembler instructions with a full memory state of the current execution. We do not have anymore the **meta-data** of the executable file.

## What we mean by "Binary Programs"? Université BORDEAUX

**Abstract Model**: All **unnecessary information** for the analysis have been removed. Only **necessary information** remains.

**Source Code**: Keep track of high-level information about the program such as variables, types, functions. But also, variable and function names, and pragmas or code decorations.

**Bytecode**: May vary depending on the bytecode considered, but keep track of few high-level information about the program such as **types** and **functions**. But, programs are usually **unstructured**.

**Binary File**: Only keep track of the **instructions** in an **unstructured way** (no forloop, no clear argument passing in procedures, ...). **No type**, **no naming**. But, the binary file may enclose **meta-data** that might be helpful (symbols, debug, ...).

**Memory Dump**: Pure assembler instructions with a full memory state of the current execution. We do not have anymore the **meta-data** of the executable file.

## What we mean by "Binary Programs"? Université BORDEAUX

**Abstract Model**: All **unnecessary information** for the analysis have been removed. Only **necessary information** remains.

**Source Code**: Keep track of high-level information about the program such as variables, types, functions. But also, variable and function names, and pragmas or code decorations.

**Bytecode**: May vary depending on the bytecode considered, but keep track of few high-level information about the program such as **types** and **functions**. But, programs are usually **unstructured**.

**Binary File**: Only keep track of the **instructions** in an **unstructured way** (no forloop, no clear argument passing in procedures, ...). No type, no naming. But, the binary file may enclose **meta-data** that might be helpful (symbols, debug, ...).

**Memory Dump**: Pure assembler instructions with a full memory state of the current execution. We do not have anymore the **meta-data** of the executable file.

#### Binary code is the closest format of what will be executed!



- Loader: Open the input file, parse the meta-data enclosed in the binary file and extract the code to be mapped in memory.
- **Decoder**: Given a **sequence of bytes** at an address in memory, translate it into an **intermediate representation** which will be analyzed afterward.
- **Disassembler**: Combination of a **decoder** and a **strategy** to browse through the memory in order to recover all the control-flow of the program.
- **Decompiler**: Translate the assembly code into a high-level language with variables, types, functions and more (modules, objects, classes, ...).
- Verificator: Take the high-level representation of the program and check it against formally specified properties.



- Loader: Open the input file, parse the meta-data enclosed in the binary file and extract the code to be mapped in memory.
- **Decoder**: Given a **sequence of bytes** at an address in memory, translate it into an **intermediate representation** which will be analyzed afterward.
- **Disassembler**: Combination of a **decoder** and a **strategy** to browse through the memory in order to recover all the control-flow of the program.
- **Decompiler**: Translate the assembly code into a high-level language with variables, types, functions and more (modules, objects, classes, ...).
- Verificator: Take the high-level representation of the program and check it against formally specified properties.

- Trustable reconstruction of the program control-flow;
- "As much as we can" automation of recovery of the control-flow;
- Scaling the analysis from small to big binary software;
- Performing automatic and correct, but partial, decompilation;
- Verification of few accessibility properties on real binary programs;

<sup>de</sup> BORDFAIIX

- Trustable reconstruction of the program control-flow;
- "As much as we can" automation of recovery of the control-flow;
- Scaling the analysis from small to big binary software;
- Performing automatic and correct, but partial, decompilation;
- Verification of few accessibility properties on real binary programs;

It does not seems to be a lot, but it is already quite tricky! université

<sup>de</sup> BORDFAIIX

université BORDEAUX

#### Introducing to Binary Code Analysis

#### 2 Why Is Binary Analysis Special?

- Unstructured Programming
- Architectural Model

#### 3 Low-level Programs Formal Model

- 4 Control-flow Recovery
- **5** Current and Future Trends

## université BORDEAUX

#### No Advanced Programming Constructs and Types

- No variable (only registers and memory accesses)
- No advanced types (only: Value, Pointer or Instructions);
- No advanced control-flow constructs (if-then-else, for, while, ...);

#### **Jump-based Programming**

- Static Jumps: jmp 0x12345678
- Dynamic Jumps: jmp \*%eax

#### **No Function Facilities**

- No Function Type or Definition;
- No Argument Passing Facilities;
- No Procedural Context Facilities;

#### Harvard Architecture

- First implemented in the Mark I (1944).
- Keep program and data separated.
- Allows to fetch data and instructions in the same time.

#### Princeton Architecture (Von Neumann)

- First implemented in the ENIAC (1946).
- Allows self-modifying code and entanglement of program and data.



université





- Princeton Ar Low-level programming





- Introducing to Binary Code Analysis
- 2 Why Is Binary Analysis Special?
- **3** Low-level Programs Formal Model
- Control-flow Recovery
- 5 Current and Future Trends

- Semantics of low-level programs differ drastically from the usual models;
- Real execution models are optimized a lot which make them difficult to handle;
- A simpler model with the same expressivity make it easier to understand;
- A formalization is necessary to start thinking about proofs;

# Memory

- $\mathbb{D} \subseteq \mathbb{N}$ : A discrete numerical domain;
- $\mathbb{A} = \mathbb{D}$ : Memory addresses (part of the numerical domain);
- $\mathbb{M} : \mathbb{A} \mapsto \mathbb{D}$ : The set of all possible valuations of the memory;
- Notation:  $m \in \mathbb{M}$ , m(addr) = val.

# Partially Initialized Memory

- $\mathbb{M}|_A : \mathbb{A} \mapsto \mathbb{D} \cup \{\bot\}$ : The set of all partial valuations of  $\mathbb{M}$ , with  $A \subseteq \mathbb{A}$  the initialized addresses such that  $\forall a \in \mathbb{A} \setminus A$ ,  $m(a) = \bot$ .
- Notation: If  $m \in \mathbb{M}|_A$ , then  $\mathbb{M}(m)$  denotes the set of all the fully initialized memories that can be spawned with m as generator.

# Register(s)

 $\bullet \mbox{ pc} \in \mathbb{A}:$  The program counter (the only register of the model);

université

<sup>de</sup> BORDFAIIX

# Assembly Language

#### Instructions

- I: A (finite) set of instructions;
- 'load value, addr': Load the evaluation of 'value' at 'addr' in memory;
- 'branch cond, addr': Jump to 'addr' if the expression 'cond' is zero;
- 'halt': Stop program execution;

#### Expressions

Expressions are usual arithmetics (e.g. '10\*(5-7)/3') with:

•  $[addr] \in \mathbb{D}$ : Access to the content of the address 'addr'  $\in \mathbb{A}$ ;

## **Operational Semantics**

- $\mathbb{I}: \mathbb{M} \times \mathbb{A} \mapsto \mathbb{M} \times \mathbb{A}$  where  $i \in \mathbb{I}$ , i(m, pc) = (m', pc');
- [load value, addr]] = ([addr]:=value, pc':=pc+1)
- [[branch cond, addr]] =
   ([0]:=[0], if cond==0 then pc':=addr else pc':=pc+1)
- [[halt]] = ([0]:=[0], pc':=pc)

## System Calls (optional)

- syscall read addr: Get an input (keyboard) and store it into 'addr';
- syscall write value: Write 'value' on the output (screen).

université

Université BORDEAUX

# **Decoding Instructions**

- I: A set of instructions as described before;
- $\delta:\mathbb{D}\mapsto\mathbb{I}:$  A decoding function to map a value to an instruction.

# Low-Level Program

- A program  $P = (m_{init}, pc_0, \delta)$ , is given by:
  - An initial, partially initialized, memory  $m_{init} \in \mathbb{M}|_{\mathcal{A}}$  (with  $\mathcal{A} \subseteq \mathbb{A}$ ),
  - An initial program counter  $\mathtt{pc}_0\in\mathbb{A},$
  - And a decoding function  $\delta : \mathbb{D} \mapsto \mathbb{I}$ .

# Valid Run

$$(m_0, \mathrm{pc}_0) \xrightarrow{i_0(m_0, \mathrm{pc}_0)} (m_1, \mathrm{pc}_1) \xrightarrow{i_1(m_1, \mathrm{pc}_1)} \dots \xrightarrow{i_k(m_k, \mathrm{pc}_k)} (m_{k+1}, \mathrm{pc}_{k+1}) \dots$$

Where  $m_0 \in \mathbb{M}(m_{init})$  and  $\forall p \ge 0$ ,  $i_p = \delta(m_p, pc_p)$  and  $(m_{p+1}, pc_{p+1}) = i_p(m_p, pc_p)$ .

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta:$  We already applied it to the memory when needed.

Addr	Initial Content
0x0	$\perp$
0x1	$\perp$
0x2	syscall read O
0x3	load [0], 1
0x4	load [0]*[1], 1
0x5	load [0]-1, 0
0x6	branch [0]!=0, 4
0x7	branch [1]!=0, 9
0x8	load 1, [1]
0x9	syscall write [1]
0xa	halt

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta :$  We already applied it to the memory when needed.

Addr	Initial Content
0x0	$\perp$
0x1	$\perp$
0x2	syscall read O
0x3	load [0], 1
0x4	load [0]*[1], 1
0x5	load [0]-1, 0
0x6	branch [0]!=0, 4
0x7	branch [1]!=0, 9
0x8	load 1, [1]
0x9	syscall write [1]
0xa	halt

;; counter (var)

;; accumulator (var)

université

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta :$  We already applied it to the memory when needed.

Addr	Initial Content
0x0	$\perp$
0x1	$\perp$
0x2	syscall read O
0x3	load [0], 1
0x4	load [0]*[1], 1
0x5	load [0]-1, 0
0x6	branch [0]!=0, 4
0x7	branch [1]!=0, 9
0x8	load 1, [1]
0x9	syscall write [1]
0xa	halt

- ;; counter (var)
- ;; accumulator (var)
- ;; get initial value

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta$ : We already applied it to the memory when needed.

Addr	Initial Content
0x0	$\perp$
0x1	$\perp$
0x2	syscall read O
0x3	load [0], 1
0x4	load [0]*[1], 1
0x5	load [0]-1, 0
0x6	branch [0]!=0, 4
0x7	branch [1]!=0, 9
0x8	load 1, [1]
0x9	syscall write [1]
0xa	halt

- ;; counter (var)
  ;; accumulator (var)
  ;; get initial value
  ;; initialize accumulator

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta$ : We already applied it to the memory when needed.

Addr	Initial Content
0x0	$\perp$
0x1	$\perp$
0x2	syscall read O
0x3	load [0], 1
0x4	load [0]*[1], 1
0x5	load [0]-1, 0
0x6	branch [0]!=0, 4
0x7	branch [1]!=0, 9
0x8	load 1, [1]
0x9	syscall write [1]
0xa	halt

- ;; counter (var)
  ;; accumulator (var)
- ;; get initial value
- ;; initialize accumulator
- ;; compute next step

- $m_0$  as below:
- pc<sub>0</sub> = 2;
- $\delta$ : We already applied it to the memory when needed.

Addr	Initial Content
0x0	$\perp$
0x1	$\perp$
0x2	syscall read O
0x3	load [0], 1
0x4	load [0]*[1], 1
0x5	load [0]-1, 0
0x6	branch [0]!=0, 4
0x7	branch [1]!=0, 9
0x8	load 1, [1]
0x9	syscall write [1]
0xa	halt

- ;; counter (var)
  ;; accumulator (var)
  ;; get initial value
  - ;; initialize accumulator
  - ;; compute next step
  - ;; decrement counter

- $m_0$  as below:
- pc<sub>0</sub> = 2;
- $\delta$ : We already applied it to the memory when needed.

Addr	Initial Content
0x0	$\perp$
0x1	$\perp$
0x2	syscall read O
0x3	load [0], 1
0x4	load [0]*[1], 1
0x5	load [0]-1, 0
0x6	branch [0]!=0, 4
0x7	branch [1]!=0, 9
0x8	load 1, [1]
0x9	syscall write [1]
0xa	halt

- ;; counter (var)
  ;; accumulator (var)
  ;; get initial value
  ;; initialize accumulator
  - ;; compute next step
  - ;; decrement counter
    ;; loop if counter is not zero

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta :$  We already applied it to the memory when needed.

Addr	Initial Content
0x0	$\perp$
0x1	$\perp$
0x2	syscall read 0
0x3	load [0], 1
0x4	load [0]*[1], 1
0x5	load [0]-1, 0
0x6	branch [0]!=0, 4
0x7	branch [1]!=0, 9
0x8	load 1, [1]
0x9	syscall write [1]
0xa	halt

- ;; counter (var)
- ;; accumulator (var)
- ;; get initial value
- ;; initialize accumulator
- ; compute next step
- ;; decrement counter
- ;; loop if counter is not zero
- ;; check if result is not zero
- ;; if result was zero, set result to 1

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta :$  We already applied it to the memory when needed.

Addr	Initial Content
0x0	$\perp$
0x1	$\perp$
0x2	syscall read O
0x3	load [0], 1
0x4	load [0]*[1], 1
0x5	load [0]-1, 0
0x6	branch [0]!=0, 4
0x7	branch [1]!=0, 9
0x8	load 1, [1]
0x9	syscall write [1]
0xa	halt

- ;; counter (var)
- ;; accumulator (var)
- ;; get initial value
- ;; initialize accumulator
- ; compute next step
- ;; decrement counter
- ;; loop if counter is not zero
- ;; check if result is not zero
- ;; if result was zero, set result to 1
- ;; output result

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta :$  We already applied it to the memory when needed.

Addr	Initial Content	
0x0	$\perp$	;; counter (var)
0x1	$\perp$	;; accumulator (var)
0x2	syscall read 0	;; get initial value
0x3	load [0], 1	;; initialize accumulator
0x4	load [0]*[1], 1	;; compute next step
0x5	load [0]-1, 0	;; decrement counter
0x6	branch [0]!=0, 4	;; loop if counter is not zero
0x7	branch [1]!=0, 9	;; check if result is not zero
0x8	load 1, [1]	;; if result was zero, set result to 1
0x9	syscall write [1]	;; output result
0xa	halt	;; halt program

Université BORDEAUX

- m<sub>0</sub> as below;
- $pc_0 = 1;$
- $\delta :$  We already applied it to the memory when needed.

Addr	Initial Content
0x0	$\perp$
0x1	syscall read O
0x2	branch 0<[1]<4, [1]*2+2
0x3	branch 0==0, 1
0x4	syscall write 10
0x5	halt
0x6	syscall write 42
0x7	halt
0x8	syscall write 1001
0x9	halt

- m<sub>0</sub> as below;
- $pc_0 = 1;$
- $\delta :$  We already applied it to the memory when needed.

Addr	Initial Content	
0x0	$\perp$	;
0x1	syscall read 0	
0x2	branch 0<[1]<4, [1]*2+2	
0x3	branch 0==0, 1	
0x4	syscall write 10	
0x5	halt	
0x6	syscall write 42	
0x7	halt	
0x8	syscall write 1001	
0x9	halt	

;; input (var)

université

- *m*<sub>0</sub> as below;
- $pc_0 = 1;$
- $\delta :$  We already applied it to the memory when needed.

Addr	Initial Content	
0x0		;; input (var)
0x1	syscall read O	;; get initial value
0x2	branch 0<[1]<4, [1]*2+2	
0x3	branch 0==0, 1	
0x4	syscall write 10	
0x5	halt	
0x6	syscall write 42	
0x7	halt	
0x8	syscall write 1001	
0x9	halt	

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 1;
- $\delta$ : We already applied it to the memory when needed.

Addr	Initial Content	
0x0	$\perp$	;; input (v
0x1	syscall read O	;; get init
0x2	branch 0<[1]<4, [1]*2+2	;; dynamic
0x3	branch 0==0, 1	
0x4	syscall write 10	
0x5	halt	
0x6	syscall write 42	
0x7	halt	
0x8	syscall write 1001	
0x9	halt	

- ar)
- ial value
- jump

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 1;
- $\delta :$  We already applied it to the memory when needed.

Addr	Initial Content	
0x0	$\perp$	;; input (var)
0x1	syscall read O	;; get initial value
0x2	branch 0<[1]<4, [1]*2+2	;; dynamic jump
0x3	branch $0==0, 1$	;; loop on wrong choice
0x4	syscall write 10	
0x5	halt	
0x6	syscall write 42	
0x7	halt	
0x8	syscall write 1001	
0x9	halt	

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 1;
- $\delta :$  We already applied it to the memory when needed.

Addr	Initial Content	
0x0	$\perp$	;; input (var)
0x1	syscall read O	;; get initial value
0x2	branch 0<[1]<4, [1]*2+2	;; dynamic jump
0x3	branch 0==0, 1	;; loop on wrong choice
0x4	syscall write 10	;; output 10 on 1
0x5	halt	
0x6	syscall write 42	;; output 42 on 2
0x7	halt	
0x8	syscall write 1001	;; output 1001 on 3
0x9	halt	]
- *m*<sup>0</sup> as below;
- pc<sub>0</sub> = 2;
- $\delta$ : We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch [0]!=0, 4}$
  - 1  $\mapsto$  branch 0==0, 8

Addr	Initial Content
0x0	$\perp$
0x1	0
0x2	syscall read O
0x3	load [1], 6
0x4	load [0], 1
0x5	load [0]-1, [0]
0x6	load [1], 0
0x7	branch 0==0, 3
0x8	halt

- *m*<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta$ : We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch [0]!=0, 4}$
  - 1  $\mapsto$  branch 0==0, 8

Addr	Initial Content
0x0	$\perp$
0x1	0
0x2	syscall read O
0x3	load [1], 6
0x4	load [0], 1
0x5	load [0]-1, [0]
0x6	load [1], 0
0x7	branch 0==0, 3
0x8	halt

- ;; input (var)
- ;; initialized data

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta$ : We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch [0]!=0, 4}$
  - 1  $\mapsto$  branch 0==0, 8

	Addr	Initial Content
	0x0	$\perp$
	0x1	0
⇒	0x2	syscall read O
	0x3	load [1], 6
	0x4	load [0], 1
	0x5	load [0]-1, [0]
	0x6	load [1], 0
	0x7	branch 0==0, 3
	0x8	halt

- ;; input (var)
- ;; initialized data
- ;; get initial value

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta$ : We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch [0]!=0, 4}$
  - 1  $\mapsto$  branch 0==0, 8

	Addr	Initial Content
	0x0	n
	0x1	0
⇒	0x2	syscall read O
	0x3	load [1], 6
	0x4	load [0], 1
	0x5	load [0]-1, [0]
	0x6	load [1], 0
	0x7	branch 0==0, 3
	0x8	halt

- ;; input (var)
- ;; initialized data
- ;; get initial value

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta$ : We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch [0]!=0, 4}$
  - 1  $\mapsto$  branch 0==0, 8

	Addr	Initial Content
	0x0	n
	0x1	0
	0x2	syscall read O
$\Rightarrow$	0x3	load [1], 6
	0x4	load [0], 1
	0x5	load [0]-1, [0]
	0x6	load [1], 0
	0x7	branch 0==0, 3
	0x8	halt

- ;; input (var)
- ;; initialized data
- ;; get initial value
- ;; rewrite code ahead

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta$ : We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch [0]!=0, 4}$
  - 1  $\mapsto$  branch 0==0, 8

	Addr	Initial Content
	0x0	n
	0x1	0
	0x2	syscall read 0
$\Rightarrow$	0x3	load [1], 6
	0x4	load [0], 1
	0x5	load [0]-1, [0]
	0x6	branch [0]!=0, 4
	0x7	branch 0==0, 3
	0x8	halt

- ;; input (var)
- ;; initialized data
- ;; get initial value
- ;; rewrite code ahead

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;

=

- $\delta$ : We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch [0]!=0, 4}$
  - 1  $\mapsto$  branch 0==0, 8

	Addr	Initial Content
	0x0	n
	0x1	0
	0x2	syscall read O
	0x3	load [1], 6
>	0x4	load [0], 1
	0x5	load [0]-1, [0]
	0x6	branch [0]!=0, 4
	0x7	branch 0==0, 3
	0x8	halt

;; input (var)
;; initialized data
;; get initial value
;; rewrite code ahead
:: overwrite [1] with [0]

université

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;

=

- $\delta$ : We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch [0]!=0, 4}$
  - 1  $\mapsto$  branch 0==0, 8

Addr	Initial Content
0x0	n
0x1	n
0x2	syscall read O
0x3	load [1], 6
0x4	load [0], 1
0x5	load [0]-1, [0]
0x6	branch [0]!=0, 4
0x7	branch 0==0, 3
0x8	halt

;; input (var) ;; initialized data ;; get initial value ;; rewrite code ahead :: overwrite [1] with [0]

université

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta$ : We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch [0]!=0, 4}$
  - 1  $\mapsto$  branch 0==0, 8

	Addr	Initial Content
	0x0	n
	0x1	n
	0x2	syscall read O
	0x3	load [1], 6
	0x4	load [0], 1
$\Rightarrow$	0x5	load [0]-1, [0]
	0x6	branch [0]!=0, 4
	0x7	branch 0==0, 3
	0x8	halt

- ;; input (var) ;; initialized data ;; get initial value
- ;; rewrite code ahead
- ;; overwrite [1] with [0]
- ;; decrement [0]

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta$ : We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch [0]!=0, 4}$
  - 1  $\mapsto$  branch 0==0, 8

	Addr	Initial Content
	0x0	n-1
	0x1	n
	0x2	syscall read O
	0x3	load [1], 6
	0x4	load [0], 1
$\Rightarrow$	0x5	load [0]-1, [0]
	0x6	branch [0]!=0, 4
	0x7	branch 0==0, 3
	0x8	halt

- ;; input (var)
  ;; initialized data
- ;; get initial value
- ;; rewrite code ahead
- ;; overwrite [1] with [0]
- ;; decrement [0]

Self-modifying code

- $pc_0 = 2;$
- $\delta$ : We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch } [0] !=0, 4$
  - 1  $\mapsto$  branch 0==0, 8

Γ	Addr	Initial Content
Γ	0x0	n-1
Γ	0x1	n
Γ	0x2	syscall read 0
	0x3	load [1], 6
	0x4	load [0], 1
	0x5	load [0]-1, [0]
> [	0x6	branch [0]!=0, 4
	0x7	branch 0==0, 3
ſ	0x8	halt

- ;; input (var)
- ; initialized data
- ;; get initial value
- ;; rewrite code ahead
- ;; overwrite [1] with [0]
- ; decrement [0]
- ; if not zero loop to 4

université

• m<sub>0</sub> as below;

Self-modifying code

- $\delta$ : We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch } [0] !=0, 4$
  - $1 \mapsto \text{branch } 0==0, 8$

	Addr	Initial Content	
	0x0	0	;; input (var)
	0x1	1	;; initialized data
	0x2	syscall read O	;; get initial value
	0x3	load [1], 6	;; rewrite code ahead
	0x4	load [0], 1	;; overwrite [1] with [0]
	0x5	load [0]-1, [0]	;; decrement [0]
	0x6	branch [0]!=0, 4	;; if not zero loop to 4
$\Rightarrow$	0x7	branch 0==0, 3	;; jump to 3
	0x8	halt	

université

• m<sub>0</sub> as below;

Self-modifying code

- $\delta :$  We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch [0]!=0, 4}$
  - 1  $\mapsto$  branch 0==0, 8

	Addr	Initial Content	
	0x0	0	;; input (var)
	0x1	1	;; initialized data
	0x2	syscall read O	;; get initial value
$\Rightarrow$	0x3	load [1], 6	;; rewrite code ahead
	0x4	load [0], 1	;; overwrite [1] with [0]
	0x5	load [0]-1, [0]	;; decrement [0]
	0x6	branch [0]!=0, 4	;; if not zero loop to 4
	0x7	branch 0==0, 3	;; jump to 3
	0x8	halt	

université

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta$ : We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch [0]!=0, 4}$
  - 1  $\mapsto$  branch 0==0, 8

	Addr	Initial Content	
	0x0	0	;; input (var)
	0x1	1	;; initialized data
$\Rightarrow$	0x2	syscall read O	;; get initial value
	0x3	load [1], 6	;; rewrite code ahead
	0x4	load [0], 1	;; overwrite [1] with
	0x5	load [0]-1, [0]	;; decrement [0]
	0x6	branch 0==0, 8	;; jump to 8
	0x7	branch 0==0, 3	;; jump to 3
	0x8	halt	

[0]

université

- m<sub>0</sub> as below;
- $pc_0 = 2;$
- $\delta$ : We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch } [0] !=0, 4$
  - $1 \mapsto \text{branch } 0==0, 8$

	Addr	Initial Content	
	0x0	0	;; input (var)
	0x1	1	;; initialized data
	0x2	syscall read O	;; get initial value
	0x3	load [1], 6	;; rewrite code ahead
$\Rightarrow$	0x4	load [0], 1	;; overwrite [1] with
	0x5	load [0]-1, [0]	;; decrement [0]
	0x6	branch 0==0, 8	;; jump to 8
	0x7	branch 0==0, 3	;; jump to 3
	0x8	halt	
	0x8	nait	

with [0]

université

- m<sub>0</sub> as below;
- $pc_0 = 2;$
- $\delta$ : We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch } [0] !=0, 4$
  - $1 \mapsto \text{branch } 0==0, 8$

[	Addr	Initial Content	
	0x0	0	;; input (var)
	0x1	0	;; initialized data
⇒	0x2	syscall read O	;; get initial value
	0x3	load [1], 6	;; rewrite code ahead
	0x4	load [0], 1	;; overwrite [1] with [0]
	0x5	load [0]-1, [0]	;; decrement [0]
	0x6	branch 0==0, 8	;; jump to 8
	0x7	branch 0==0, 3	;; jump to 3
	0x8	halt	

- *m*<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta$ : We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch [0]!=0, 4}$
  - 1  $\mapsto$  branch 0==0, 8

	Addr	Initial Content	
	0x0	0	;; inp
	0x1	0	;; ini
	0x2	syscall read O	;; get
	0x3	load [1], 6	;; rew
	0x4	load [0], 1	;; ove
$\Rightarrow$	0x5	load [0]-1, [0]	;; dec
	0x6	branch 0==0, 8	;; jum
	0x7	branch 0==0, 3	;; jum
	0x8	halt	

; input (var) ; initialized data ; get initial value ; rewrite code ahead ; overwrite [1] with [0] ; decrement [0] ; jump to 8 ; jump to 3

université

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta$ : We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch [0]!=0, 4}$
  - 1  $\mapsto$  branch 0==0, 8

	Addr	Initial Content	
	0x0	-1	;; inpu
	0x1	0	;; init
	0x2	syscall read O	;; get
	0x3	load [1], 6	;; rewr
	0x4	load [0], 1	;; over
$\Rightarrow$	0x5	load [0]-1, [0]	;; decr
	0x6	branch 0==0, 8	;; jump
	0x7	branch 0==0, 3	;; jump
	0x8	halt	

; input (var) ; initialized data ; get initial value ; rewrite code ahead ; overwrite [1] with [0] ; decrement [0] ; jump to 8 : jump to 3

université

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta$ : We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch [0]!=0, 4}$
  - 1  $\mapsto$  branch 0==0, 8

	Addr	Initial Content	
	0x0	-1	;;
	0x1	0	;;
	0x2	syscall read 0	;;
	0x3	load [1], 6	;;
	0x4	load [0], 1	;;
	0x5	load [0]-1, [0]	;;
$\Rightarrow$	0x6	branch 0==0, 8	;;
	0x7	branch 0==0, 3	;;
	0x8	halt	

; input (var) ; initialized data ; get initial value ; rewrite code ahead ; overwrite [1] with [0] ; decrement [0] ; jump to 8 ; jump to 3

université

- m<sub>0</sub> as below;
- pc<sub>0</sub> = 2;
- $\delta$ : We already applied it to the memory when needed but here are the rest:
  - $0 \mapsto \text{branch [0]!=0, 4}$
  - 1  $\mapsto$  branch 0==0, 8

Addr	Initial Content	
0x0	-1	;;
0x1	0	;;
0x2	syscall read 0	;;
0x3	load [1], 6	;;
0x4	load [0], 1	;;
0x5	load [0]-1, [0]	;;
0x6	branch 0==0, 8	;;
0x7	branch 0==0, 3	;;
0x8	halt	

; input (var) ; initialized data ; get initial value ; rewrite code ahead ; overwrite [1] with [0] ; decrement [0] ; jump to 8 ; jump to 3

 $\Rightarrow$ 

université

A few real-world assembly languages have variable size instructions. This property is sometimes used to hide part of a program with a technique called "instruction overlapping". This property can be easily added to our model as follow.

### Instructions

- I: A (finite) set of instructions;
- 'load value, addr': Load the evaluation of 'value' at 'addr' in memory Encoded in two memory cells, first for 'load value' and second for 'address';
- 'branch cond, addr': Jump to 'addr' if the expression 'cond' is zero Encoded in two memory cells, first for 'branch cond' and second for 'address';
- 'halt': Stop program execution. Encoded in one memory cell as before;

# **Operational Semantics**

- $\mathbb{I}: \mathbb{M} \times \mathbb{A} \mapsto \mathbb{M} \times \mathbb{A}$  where  $i \in \mathbb{I}$ , i(m, pc) = (m', pc');
- [load value, addr] = ([addr]:=value, pc':=pc+2)
- [[branch cond, addr]] =
   ([0]:=[0], if cond==0 then pc':=addr else pc':=pc+2)
- [[halt]] = ([0]:=[0], pc':=pc)

université

université BORDEAUX

Introducing to Binary Code Analysis

2 Why Is Binary Analysis Special?

3 Low-level Programs Formal Model

### 4 Control-flow Recovery

- Types of Control-Flow Recovery
- Syntax-based Recovery
- Semantics-based Recovery
- Control-Flow Recovery: Summary

#### 5) Current and Future Trends

- Control-flow recovery is prior to any other work because it aims at **recovering the semantics** of the program.
- The point is to gather all the possible execution paths of the binary program for all possible inputs.
- Because of dynamic jumps and self-modifying code, the gathering of all the possible runs requires to perform data-analysis on a partial semantics of the program.
- Most of the analysis techniques work only with the complete semantics of the program (Chicken and Egg Problem).
- Thus, we need to come with new techniques...

# Correctness

- **Exact**: The disassembler outputs the exact control-flow that covers all the possible execution paths of the input program.
- **Under-approximation**: The disassembler outputs a subset of all the possible execution paths of the input program.
- **Over-approximation**: The disassembler outputs a set of execution paths that enclose the set of all possible ones.
- **Incorrect**: The disassembler outputs a set that may miss some execution paths and add some extra as well (we cannot say anything from this output).

# Techniques

Syntax-based Recovery

- Linear Sweep
- Recursive Traversal

### Semantics-based Recovery

- Concrete Execution
- Symbolic Execution

université

#### Theorem

Recovering the control-flow of a binary program is undecidable (for the general case).

### **Sketch of Proof**

- **1** Lets, first, assume that the model we presented is equivalent to a Turing machine.
- 2 Recovering all the run would requires to collect all the possible values of pc.
- Because of self-modifying code, the values pointed by the pc must also be recovered (which means that we need to track strictly more than one variable).
- Thus, we can reduce any accessibility problem for a given program to a control-flow recovery problem by adding to the original program a conditional jump to an error state. And try to see if this extra program state is in the program control-flow.
- Finally, as the accessibility problem is undecidable, the control-flow recovery problem is also undecidable for the general case.

université

- Decode the first instruction at the entrypoint and store it;
- One (syntactically) the program counter to the next instruction;
- Obecode the instruction and go to 2 if you are not out of the memory.

université

- O Decode the first instruction at the entrypoint and store it;
- Move (syntactically) the program counter to the next instruction;
- **③** Decode the instruction and go to 2 if you are not out of the memory.

# Is it adding and missing execution paths?

université

- Decode the first instruction at the entrypoint and store it;
- Move (syntactically) the program counter to the next instruction;
- Obcode the instruction and go to 2 if you are not out of the memory.

### Is it adding and missing execution paths?

Lets disassemble this piece of binary code:

0804846c:	eb04	jmp	0x804846e+4
0804846 <mark>e</mark> :	efbeadde	dd	Oxdeadbeef # Data hidden among instructions
08048472:	a16e840408	mov	eax, [0x804846e]
08048477:	83c00a	add	eax, Oxa

université

- Decode the first instruction at the entrypoint and store it;
- Move (syntactically) the program counter to the next instruction;
- Solution Decode the instruction and go to 2 if you are not out of the memory.

# Is it adding and missing execution paths?

Lets disassemble this piece of binary code:

0804846c:	eb04	jmp	0x804846e+4
0804846 <b>e</b> :	efbeadde	dd	Oxdeadbeef # Data hidden among instructions
08048472:	a16e840408	mov	eax, [0x804846e]
08048477:	83c00a	add	eax, Oxa

0804846c: eb04 jmp 0x804846e+4

université

- Decode the first instruction at the entrypoint and store it;
- Move (syntactically) the program counter to the next instruction;
- Solution Decode the instruction and go to 2 if you are not out of the memory.

# Is it adding and missing execution paths?

Lets disassemble this piece of binary code:

0804846c: 0804846e: 08048472: 08048477:	eb04 efbeadde a16e840408 83c00a	jmp dd mov add	0x804846e+4 Oxdeadbeef <b># Data hidden among instructions</b> eax, [0x804846e] eax, 0xa
0804846c:	eb04 ef	jmp out	0x804846e+4

université

- Decode the first instruction at the entrypoint and store it;
- Move (syntactically) the program counter to the next instruction;
- Solution Decode the instruction and go to 2 if you are not out of the memory.

# Is it adding and missing execution paths?

Lets disassemble this piece of binary code:

0804846c:	eb04	jmp	0x804846e+4
0804846e:	efbeadde	dd	0xdeadbeef
08048472:	a16e840408	mov	eax, [0x804846e]
08048477:	83c00a	add	eax, 0xa

0804846c:	eb04	jmp	0x804846e+4
0804846 <mark>e</mark> :	ef	out	dx, eax
0804846f:	beaddea16e	mov	esi, Ox6ea1dead

université

- Decode the first instruction at the entrypoint and store it;
- Move (syntactically) the program counter to the next instruction;
- Solution Decode the instruction and go to 2 if you are not out of the memory.

## Is it adding and missing execution paths?

Lets disassemble this piece of binary code:

0804846c:	eb04	jmp	0x804846e+4
0804846 <mark>e</mark> :	efbeadde	dd	Oxdeadbeef # Data hidden among instructions
08048472:	a16e840408	mov	eax, [0x804846e]
08048477:	83c00a	add	eax, Oxa

0804846c:	eb04	jmp	0x804846e+4
0804846 <b>e</b> :	ef	out	dx, eax
0804846f:	beaddea16e	mov	esi, Ox6ea1dead
08048474:	840408	test	[eax+ecx], al

université

- Decode the first instruction at the entrypoint and store it;
- Move (syntactically) the program counter to the next instruction;
- Solution Decode the instruction and go to 2 if you are not out of the memory.

### Is it adding and missing execution paths?

Lets disassemble this piece of binary code:

0804846c:	eb04	jmp	0x804846e+4
0804846e:	efbeadde	dd	0xdeadbeef
08048472:	a16e840408	mov	eax, [0x804846e]
08048477:	83c00a	add	eax, 0xa

0804846c:	eb04	jmp	0x804846e+4
0804846e:	ef	out	dx, eax
0804846f:	beaddea16e	mov	esi, Ox6ea1dead
08048474:	840408	test	[eax+ecx], al
08048477:	83c00a	add	eax, Oxa

université

- Decode the first instruction at the entrypoint and store it;
- Move (syntactically) the program counter to the next instruction;
- Solution Decode the instruction and go to 2 if you are not out of the memory.

## Is it adding and missing execution paths?

Lets disassemble this piece of binary code:

0804846c:	eb04	jmp	0x804846e+4
0804846 <b>e</b> :	efbeadde	dd	Oxdeadbeef # Data hidden among instructions
08048472:	a16e840408	mov	eax, [0x804846e]
08048477:	83c00a	add	eax, Oxa

0804846c:	eb04	jmp	0x804846e+4
0804846 <b>e</b> :	ef	out	dx, eax
0804846f:	beaddea16e	mov	esi, Ox6ea1dead
08048474:	840408	test	[eax+ecx], al
08048477:	83c00a	add	eax, Oxa

# Yes, it is adding and missing execution paths!

université

# Syntax-based: Linear Sweep

### Linear Sweep

- Decode the first instruction are entrypoint and store it;
- 2 Move (syntactically) the program program to the next instruction;
- Decode the instruction and go to 2 if the pre not out of the memory.

### Is it adding and missing execution paths?

Lets disassemble this piece of binary code:

0804846c:	eb04	jmp	0x804846e+4
0804846 <mark>e</mark> :	efbeadde	dd	Oxdeadbeef # Data hidden among instructions
08048472:	a16e840408	mov	eax, [0x804846e]
08048477:	83c00a	add	eax, Oxa

0804846c:	eb04	jmp	0x804846e+4
0804846 <b>e</b> :	ef	out	dx, eax
0804846f:	beaddea16e	mov	esi, Ox6ea1dead
08048474:	840408	test	[eax+ecx], al
08048477:	83c00a	add	eax, Oxa

### Yes, it is adding and missing execution paths!

université

# Syntax-based: Recursive Traversal

Introduce a partial support of one type of dynamic jump (call/ret) with almost no semantics support.

**Recursive Traversal** 

- Do linear sweep until encountering a 'call' or a 'ret';
- If this is a 'call', stack its address, jump to it and go to 1;
- If this is a 'ret', pop the last address from the stack, jump to it and go to 1.

université

BORDEAUX
Introduce a partial support of one type of dynamic jump (call/ret) with almost no semantics support.

**Recursive Traversal** 

- Do linear sweep until encountering a 'call' or a 'ret';
- If this is a 'call', stack its address, jump to it and go to 1;
- If this is a 'ret', pop the last address from the stack, jump to it and go to 1.

## What does it add to linear sweep?

université

BORDEAUX

Introduce a partial support of one type of dynamic jump (call/ret) with almost no semantics support.

**Recursive Traversal** 

- Do linear sweep until encountering a 'call' or a 'ret';
- If this is a 'call', stack its address, jump to it and go to 1;
- If this is a 'ret', pop the last address from the stack, jump to it and go to 1.

## What does it add to linear sweep?

Lets disassemble this piece of binary code:

```
      0804846c:
      e882feffff call 0x08048c00
      08048c00:
      83c00010 add eax, 0x1000

      08048471:
      a16e840408 mov eax, [0x804846e]
      08048c03:
      c3
      ret

      08048476:
      83c00a
      add eax, 0xa
      ...
      ...
```

université

Introduce a partial support of one type of dynamic jump (call/ret) with almost no semantics support.

**Recursive Traversal** 

- Do linear sweep until encountering a 'call' or a 'ret';
- If this is a 'call', stack its address, jump to it and go to 1;
- If this is a 'ret', pop the last address from the stack, jump to it and go to 1.

## What does it add to linear sweep?

Lets disassemble this piece of binary code:

```
      0804846c:
      e882feffff
      call
      0x08048c00
      08048c00:
      83c00010
      add
      eax, 0x1000

      08048471:
      a16e840408
      mov
      eax, [0x804846e]
      08048c03:
      c3
      ret

      08048476:
      83c00a
      add
      eax, 0xa
      ...
```

0804846c: e882feffff call 0x08048c00

université

Introduce a partial support of one type of dynamic jump (call/ret) with almost no semantics support.

**Recursive Traversal** 

- Do linear sweep until encountering a 'call' or a 'ret';
- If this is a 'call', stack its address, jump to it and go to 1;
- If this is a 'ret', pop the last address from the stack, jump to it and go to 1.

## What does it add to linear sweep?

Lets disassemble this piece of binary code:

```
      0804846c:
      e882feffff call 0x08048c00
      08048c00:
      83c00010 add eax, 0x1000

      08048471:
      a16e840408 mov eax, [0x804846e]
      08048c03:
      c3
      ret

      08048476:
      83c00a
      add eax, 0xa
      ...
      ret
```

0804846c: e882feffff call 0x08048c00 08048c00: 83c00010 add eax, 0x1000 université

Introduce a partial support of one type of dynamic jump (call/ret) with almost no semantics support.

**Recursive Traversal** 

- Do linear sweep until encountering a 'call' or a 'ret';
- If this is a 'call', stack its address, jump to it and go to 1;
- If this is a 'ret', pop the last address from the stack, jump to it and go to 1.

## What does it add to linear sweep?

Lets disassemble this piece of binary code:

```
      0804846c:
      e882feffff
      call
      0x08048c00
      08048c00:
      83c00010
      add
      eax, 0x1000

      08048471:
      a16e840408
      mov
      eax, [0x804846e]
      08048c03:
      c3
      ret

      08048476:
      83c00a
      add
      eax, 0xa
      ...
```

0804846c: e882feffff call 0x08048c00 08048c00: 83c00010 add eax, 0x1000 08048c03: c3 ret université

Introduce a partial support of one type of dynamic jump (call/ret) with almost no semantics support.

**Recursive Traversal** 

- Do linear sweep until encountering a 'call' or a 'ret';
- If this is a 'call', stack its address, jump to it and go to 1;
- If this is a 'ret', pop the last address from the stack, jump to it and go to 1.

## What does it add to linear sweep?

Lets disassemble this piece of binary code:

```
      0804846c:
      e882feffff
      call
      0x08048c00
      08048c00:
      83c00010
      add
      eax, 0x1000

      08048471:
      a16e840408
      mov
      eax,
      [0x804846e]
      08048c03:
      c3
      ret

      08048476:
      83c00a
      add
      eax,
      0xa
      ...
```

0804846c: e882feffff call 0x08048c00 08048c00: 83c00010 add eax, 0x1000 08048c03: c3 ret 08048471: a16e840408 mov eax, [0x804846e] université

Introduce a partial support of one type of dynamic jump (call/ret) with almost no semantics support.

**Recursive Traversal** 

- Do linear sweep until encountering a 'call' or a 'ret';
- If this is a 'call', stack its address, jump to it and go to 1;
- If this is a 'ret', pop the last address from the stack, jump to it and go to 1.

## What does it add to linear sweep?

Lets disassemble this piece of binary code:

```
      0804846c:
      e882feffff
      call
      0x08048c00
      08048c00:
      83c00010
      add
      eax, 0x1000

      08048471:
      a16e840408
      mov
      eax, [0x804846e]
      08048c03:
      c3
      ret

      08048476:
      83c00a
      add
      eax, 0xa
      ...
```

0804846c: e882feffff call 0x08048c00 08048c00: 83c00010 add eax, 0x1000 08048c03: c3 ret 08048471: a16e840408 mov eax, [0x804846e] 08048477: 83c00a add eax, 0xa université

Introduce a partial support of one type of dynamic jump (call/ret) with almost no semantics support.

**Recursive Traversal** 

- Do linear sweep until encountering a 'call' or a 'ret';
- If this is a 'call', stack its address, jump to it and go to 1;
- If this is a 'ret', pop the last address from the stack, jump to it and go to 1.

## What does it add to linear sweep?

Lets disassemble this piece of binary code:

```
      0804846c:
      e882feffff
      call
      0x08048c00
      08048c00:
      83c00010
      add
      eax,
      0x1000

      08048471:
      a16e840408
      mov
      eax,
      [0x804846e]
      08048c03:
      c3
      ret

      08048476:
      83c00a
      add
      eax,
      0xa
      ret
```

0804846c: e882feffff call 0x08048c00 08048c00: 83c00010 add eax, 0x1000 08048c03: c3 ret 08048471: a16e840408 mov eax, [0x804846e] 08048477: 83c00a add eax, 0xa ... université

Introduce a partial support of one type of dynamic jump (call/ret) with almost no semantics support.

**Recursive Traversal** 

- Do linear sweep until encountering a 'call' or a 'ret';
- If this is a 'call', stack its address, jump to it and go to 1;
- If this is a 'ret', pop the last address from the stack, jump to it and go to 1.

## What does it add to linear sweep?

Lets disassemble this piece of binary code:

```
      0804846c:
      e882feffff
      call
      0x08048c00
      08048c00:
      83c00010
      add
      eax,
      0x1000

      08048471:
      a16e840408
      mov
      eax,
      [0x804846e]
      08048c03:
      c3
      ret

      08048476:
      83c00a
      add
      eax,
      0xa
      ...
```

6 e
•

### But, it is based on linear sweep, so...

٦

université

Introduce a partial support of one type of dynamic jump (call/ret) with almost no semantics support.

**Recursive Traversal** 

- Do linear sweep until encountering a 'call' or a 'ret';
- If this is a 'call', stack its address, jurge to it and go to 1;
- If this is a 'ret', pop the last address from the stack, jump to it and go to 1.

## What does it add to linear sweep?

Lets disassemble this piece of binary code:

```
      0804846c:
      e882feffff
      call
      0x08048c00
      08048c00:
      83c00010
      add
      eax,
      0x1000

      08048471:
      a16e840408
      mov
      eax,
      [0x804846e]
      08048c03:
      c3
      ret

      08048476:
      83c00a
      add
      eax,
      0xa
      ...
```

0804846c:	e882feffff	call 0x08048c00
08048c00:	83c00010	add eax, 0x1000
08048c03:	c3	ret
08048471:	a16e840408	mov eax, [0x804846e
08048477:	83c00a	add eax, Oxa

### But, it is based on linear sweep, so...

٦

université

## What can we deduce from these examples?

Having partial knowledge of the semantics, will **always** lead to miss some behaviours and produce an incorrect control-flow.

université

## What can we deduce from these examples?

Having partial knowledge of the semantics, will **always** lead to miss some behaviours and produce an incorrect control-flow.

# To be correct, a disassembler always need to know about the semantics of all the instructions!

université

#### **Concrete Execution**

Given some chosen inputs, run the program several times and collect the traces. The collection of all the traces will give you the semantics of the program.

- Efficient and simple to settle down (by using Pin, for example).
- Quite fast for a run, even if you need to store all the traces.
- Can be automatized with random inputs (fuzzing).

## But!

- There is, almost, no hope to reach full coverage of the program.
- Random input makes it very difficult to control the time needed to reach a good coverage.

université

<sup>de</sup> BORDFAIIY

#### **Concrete Execution**

Given some chosen inputs, run the program several times and collect the traces. The collection of a the traces will give you the semantics of the program.

- Efficient and simple to settle down (by using Pin, for example).
- Quite fast for a run, even if you need to store all the traces.
- Can be automatized with radom inputs (fuzzing).
- But!
  There is, almost, no hope to reach full coverage of the program.
- Random input makes it very difficult to control the time needed to reach a good coverage.

université

# Symbolic Execution

Université Bordeaux





- line 4: (x = y)
- line 8:  $(x = y) \land (y = x + 10)$  (UNSAT)
- line 10 (path1):  $(x \neq y)$
- line 10 (path2):  $(x = y) \land (y \neq x + 10)$

#### Algorithm (James King, 1976)

Explore the program and ask the SMT-solver at each program point if the path is feasible.

## Symbolic Execution

Université Bordeaux





- line 4: (x = y)
- line 8:  $(x = y) \land (y = x + 10)$  (UNSAT)
- line 10 (path1):  $(x \neq y)$
- line 10 (path2):  $(x = y) \land (y \neq x + 10)$

#### Algorithm (James King, 1976)

Explore the program and ask the SMT-solver at each program point if the path is feasible.

Directed Automated Concrete Execution

#### **Directed Automated Concrete Execution**

- First run the program on random inputs and get a trace;
- Get each possible branching inside the previous trace and ask an SMT-solver to solve it.
- If the SMT-solver fails, generate a random input to try to reach the untouched branches.

## • Original idea (2005):

DART (Directed Automated Random Testing) by Patrice Godefroid;

## • First applied to binary analysis (2008):

Inside the OSMOSE software by CEA List.

Directed Automated Concrete Execution

#### Directed Automated Concrete Execution

- First run the program on random inputs and get a trace;
- 3 Get each possible bracking inside the previous trace and ask an SMT-solver to solve it.
- If the SMT-solver fails, generate a random input to try to reach the Original idea (2005): DART (Directed Automated Random Testing) by Patrice Codefroid;
  - Original idea (2005):

• First applied to binary analysis (2008): Inside the OSMOSE software by CEA List.

# Full Symbolic Execution on Binary Code Univers

### Algorithm

- Start at entry point;
- Symbolically execute the current instruction;
- If a dynamic jump or a test is encountered, run the SMT-solver on the conjunction of all previous paths and list possible outputs;
- If the SMT-solver output an answer, follow the satisfiable paths and go to 2;
- If the SMT-solver cannot answer, stop here.
- A few limitations and challenges:
  - Tool must be aware of the semantics of all the instructions;
  - Context of the Operating System must be simulated;
  - Under-approximation (efficiency depends upon the cleverness of SMT-solver);
  - Loops are unfolded up to a certain limit to enforce termination;
  - Detection of local context and scope helps to keep the formula small.

# Full Symbolic Execution on Binary Code

## Algorithm

- Start at entry point;
- Symbolically execute the current instruction;
- If a dynamic jump or a test is encountered, run the SMT-solver on the conjunction of all previous paths and list possible outputs;
- If the SMT-solver output an answer, follow the satisfiable paths and go to 2;
- If the SMT-solver cannot answer, stor pere.

A few limitations and challenges:

- Tool must be aware of the semantics of all the instructions;
- Context of the Operating System must be simulated;
- Under-approximation (efficiency depends upon the cleverness of SMT-solver);
- Loops are unfolded up to a certain limit to enforce termination;
- Detection of local context and scope helps to keep the formula small.

# Abstract Interpretation-Based Recovery

Using an abstract interpretation framework on the CFG recovery problem is difficult because of the 'chicken-and-egg' problem.

#### Abstract Interpretation-Based CFG Recovery

In 'An abstract interpretation-based framework for control flow reconstruction from binaries' by Johannes Kinder, Florian Zuleger, and Helmut Veith (2009).

- $\bullet~$  Use a double abstract domain: CFG  $\times~$  Data-flow analysis;
- Recovery of the CFG is part of part of the process for reaching the fix-point.
- Data-flow analysis help on the way for the fix-point.
- The abstract domain of the data-flow analysis is a parameter of the framework. It can be anything as long as it match usual hypothesis of abstract domain (Galois connection, monotonicity, ...)
- Possible domains to use: k-sets, (stridded) intervals or Value-Set Analysis.

# Abstract Interpretation-Based Recovery Université

Using an abstract interpretation framework on the CFG recovery problem is difficult because of the 'chicken-and-egg' problem.

#### Abstract Interpretation-Based CFG Recovery

In 'An abstract interpretation-based framework for control flow reconstruction from binaries' by Johannes Kinder, Florian Zuleger, and Helmut Veith (2009).

- Use a double abstract domay. CFG  $\times$  Data-flow analysis;
- Recovery of the CFG is part of part of the process for reaching the fix-point.
- Data-flow analysis help on the way for the tix-point.
- The abstract domain of the data-flow analysis is a parameter of the framework. It can be anything as long as it match usual hypothesis of abstract domain (Galois connection, monotonicity,
- Possible domains to use: k-sets, (stridded) intervals or Value-Set Analysis.

Syntax-based Disassembler	Accuracy
Linear Sweep	Incorrect
Recursive Traversal	Incorrect

• All methods are just incorrect in all cases.

Semantics-Based Disassembler	Accuracy
Concrete Execution	Under-approximation
Directed Automated Concrete Execution	Under-approximation
Full Symbolic Execution	Under-approximation
Abstract Interpretation Recovery	<b>Over-approximation</b>

- Symbolic Execution and Directed Automated Concrete Execution are of the same kind and provide under-approximation. They are useful for reverse-engineering.
- Abstract-Interpretation framework are, most of the time, too imprecise.

université



- Introducing to Binary Code Analysis
- 2 Why Is Binary Analysis Special?
- **3** Low-level Programs Formal Model
- Control-flow Recovery
- 5 Current and Future Trends

# **Current Trends**

- Multiplication of tools and frameworks (reinventing the wheel).
- Clear split between academic and industry tools (complexity of use of academic tools is currently too high).
- Still some limitations to automatically recover control-flow of everyday-life binaries and to **scale**.

## **Future Trends**

- A stable and flexible framework for binary analysis.
- Support for the main platforms (Windows, Linux, \*BSD, MacOS).
- Deal with loops and variable size inputs in a more efficient way.

université



# **Questions?**