# DEOBFUSCATION: SEMANTIC ANALYSIS TO THE RESCUE

Sébastien Bardin (CEA LIST)
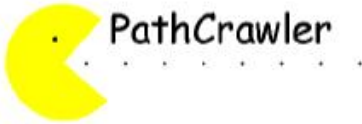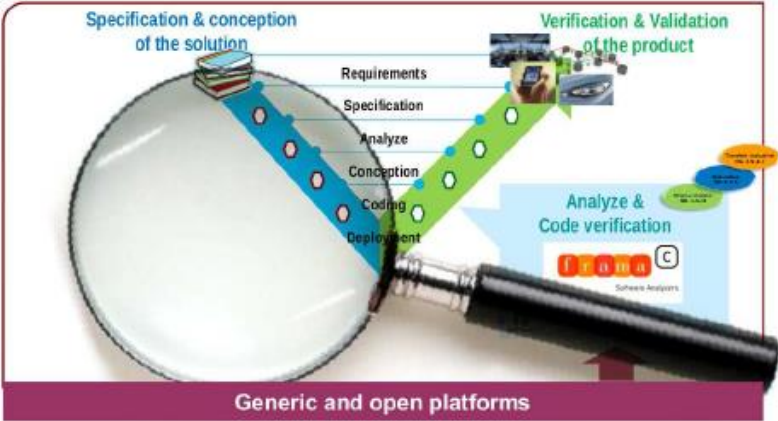Robin David (CEA LIST, QuarksLab)
Jean-Yves Marion (LORIA)

# ABOUT MY LAB @CEA     [Paris-Saclay, France]

## CEA LIST, Software Safety & Security Lab

- rigorous tools for building high-level quality software
- second part of V-cycle
- automatic software analysis
- mostly source code



Specification & conception of the solution

Verification & Validation of the product

Requirements
Specification
Analyze
Conception
Coding
Deployment

Analyze & Code verification

Generic and open platforms

frama C — Software Analyzers

GaTeL

PathCrawler

UNISIM

POLAR SSL

# IN A NUTSHELL

- **Challenge: malware *deobfuscation***

- **Standard techniques (dynamic, syntactic) not enough**

- ***Semantic methods can help*** **[obfuscation preserves semantic]**
  - Yet, need to be strongly adapted (robustness, precision, efficiency)

- ***A tour on how symbolic methods can help***
  - *Explore and discover*
  - *Prove infeasibility [S&P 2017]*
  - *Simplify* (not covered here)

# OUTLINE

- **Context**
  - *Malware comprehension*
  - *Semantic analysis*

- **The hard journey from source to binary**
  - *Explore & Discover*
  - *Prove infeasibility*

- **A few case-studies**

- **Conclusion**

# CONTEXT: MALWARE COMPREHENSION

APT: highly sophisticated attacks
- **Targeted malware**
- **Written by experts**
- Attack: 0-days
- Defense: stealth, **obfuscation**
- **Sponsored by states or mafia**

**USA elections: DNC Hack**

The day after: **malware comprehension**
- understand what has been going on
- mitigate, fix and clean
- improve defense



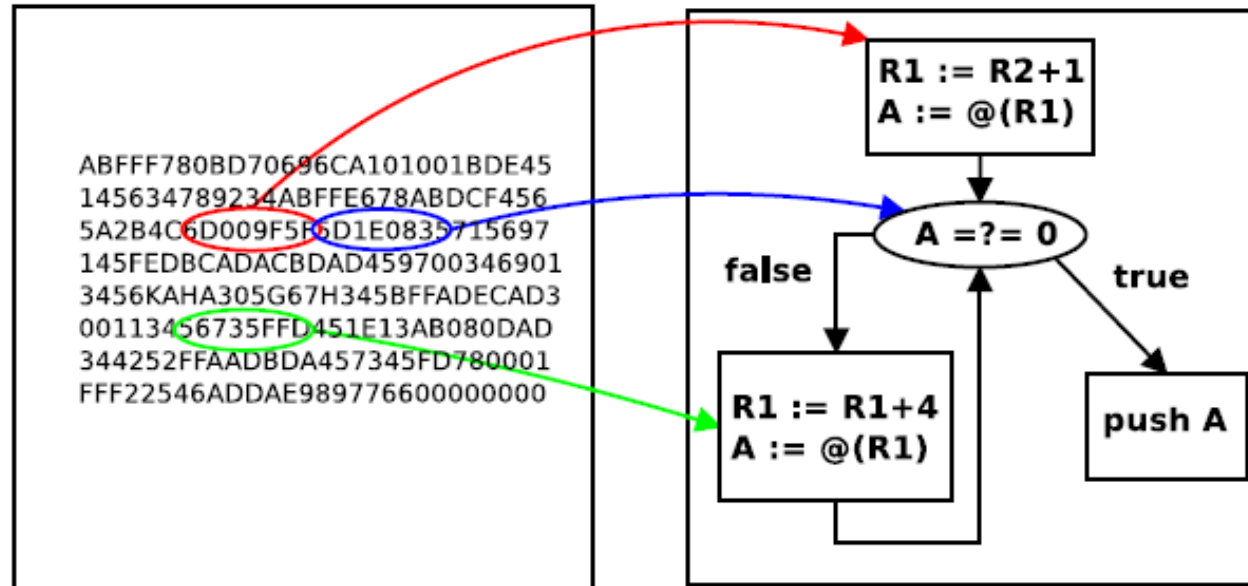**Goal: help malware comprehension**
- **Reverse of heavily obfuscated code**
- **Identify and simplify protections**

ABFFF780BD70696CA101001BDE45
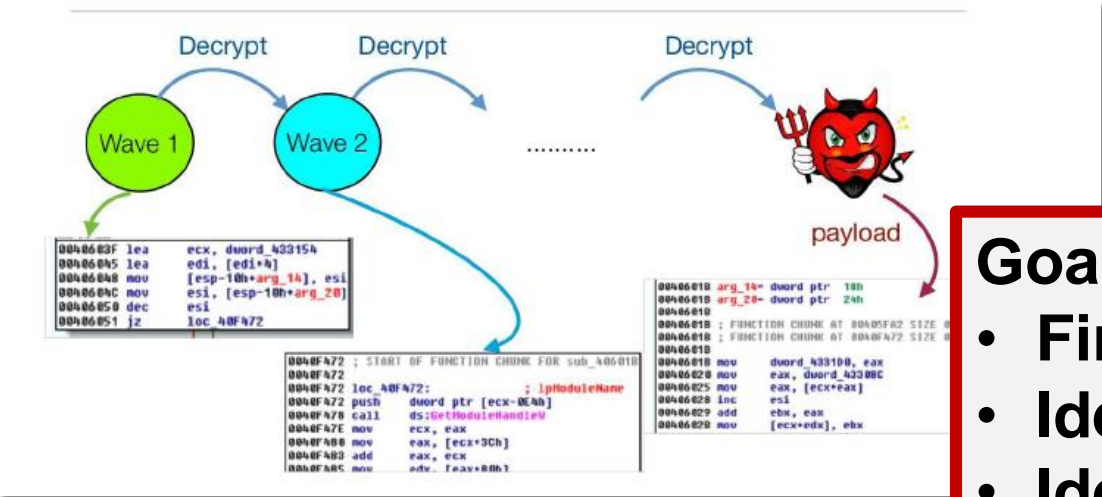145634789234ABFFE678ABDCF456
5A2B4C6D009F5F6D1E0835715697
145FEDBCADACBDAD459700346901
3456KAHA305G67H345BFFADECAD3
001134567355FFD451E13AB080DAD
344252FFAADBDA457345FD780001
FFF22546ADDAE989776600000000

R1 := R2+1
A := @(R1)

A =?= 0

false        true

R1 := R1+4
A := @(R1)

push A

**Basic reverse problem**
- **aka model recovery**
- **aka CFG recovery**

# CAN BE TRICKY!

- **code – data**
- **dynamic jumps (jmp eax)**

# REVERSE CAN BECOME A NIGHTMARE (OBFUSCATION)



**Goal: help malware comprehension**
- **Find real parts of the code**
- **Identify and simplify protections**
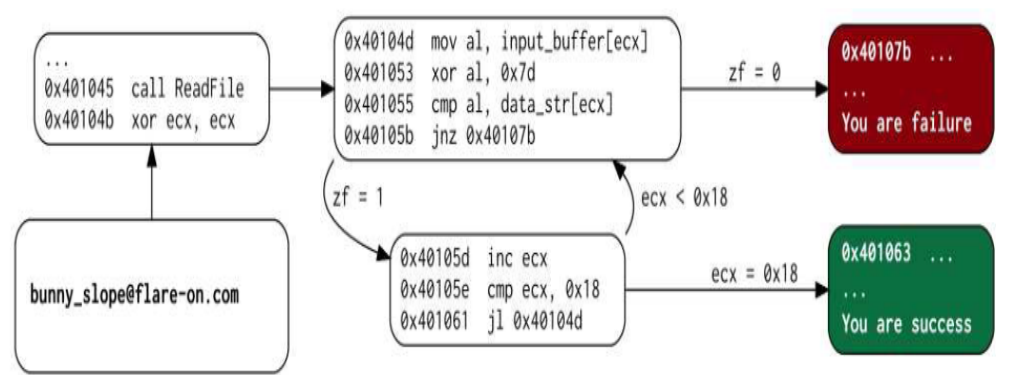- **Ideal = revert protections**

eg: $7y^2 - 1 \neq x^2$

(for any value of x, y in modular arithmetic)

```
mov   eax, ds:X
mov   ecx, ds:Y
imul  ecx, ecx
imul  ecx, 7
sub   ecx, 1
imul  eax, eax
cmp   ecx, eax
jz    <dead_addr>
```

**Obfuscation: make a code hard to reverse**
- **self-modification**
- **encryption**
- **virtualization**
- **code overlapping**
- **opaque predicates**
- **callstack tampering**
- **…**

# EXAMPLE: OPAQUE PREDICATE

**Constant-value predicates**
  (always true, always false)

- dead branch points to spurious code
- goal = waste reverser time & efforts

eg: $7y^2 - 1 \neq x^2$

(for any value of x, y in modular
arithmetic)

$\downarrow$

```
mov   eax, ds:X
mov   ecx, ds:Y
imul  ecx, ecx
imul  ecx, 7
sub   ecx, 1
imul  eax, eax
cmp   ecx, eax
jz    <dead_addr>
```
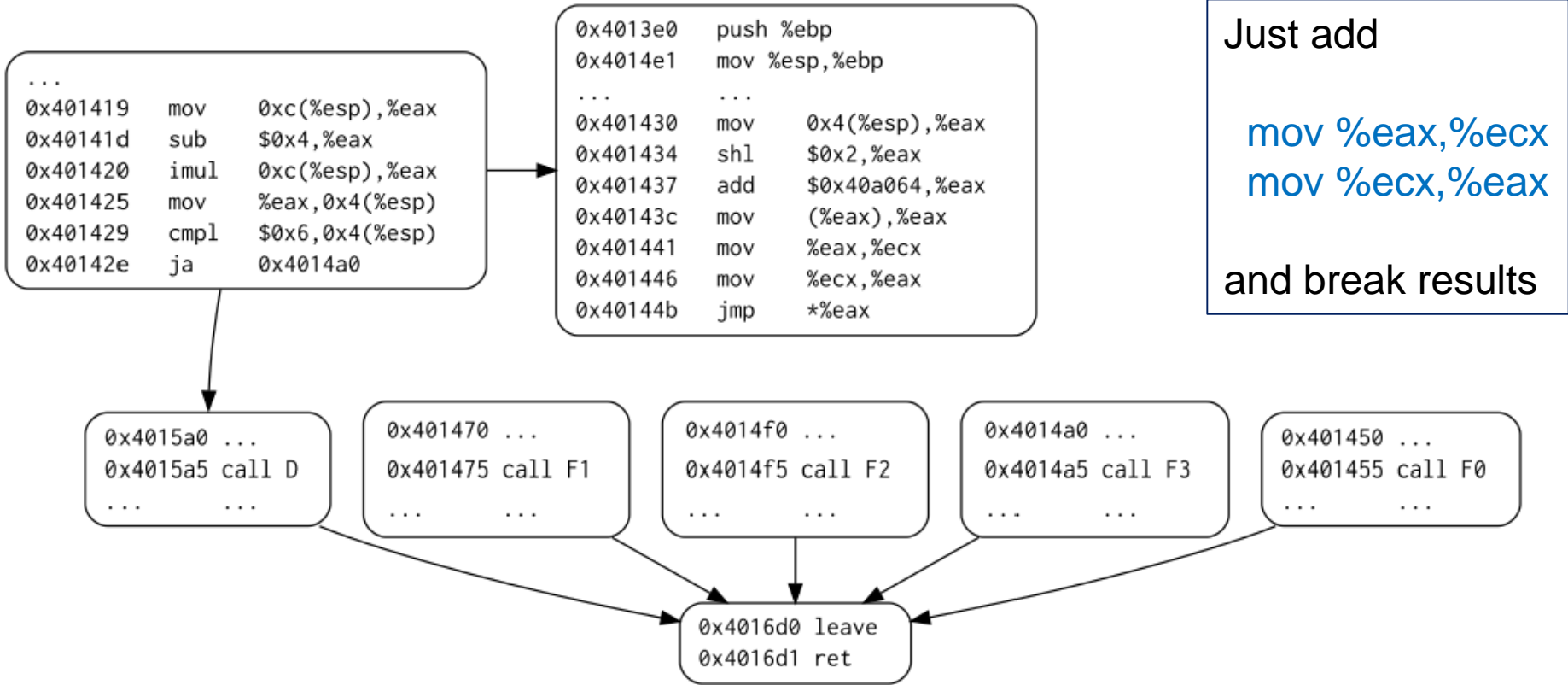
# EXAMPLE: STACK TAMPERING

**Alter the standard compilation scheme:**
   **ret do not go back to call**

- hide the real target
- return site may be spurious code

| address | instr |
|---------|-------|
| 80483d1 | call +5 |
| 80483d6 | pop edx |
| 80483d7 | add edx, 8 |
| 80483da | push edx |
| 80483db | ret |
| 80483dc | .byte{invalid} |
| 80483de | [...] |

```
0x4013e0    push %ebp
0x4014e1    mov %esp,%ebp
...         ...
0x401430    mov    0x4(%esp),%eax
0x401434    shl    $0x2,%eax
0x401437    add    $0x40a064,%eax
0x40143c    mov    (%eax),%eax
0x401441    mov    %eax,%ecx
0x401446    mov    %ecx,%eax
0x40144b    jmp    *%eax
```

```
...
0x401419    mov    0xc(%esp),%eax
0x40141d    sub    $0x4,%eax
0x401420    imul   0xc(%esp),%eax
0x401425    mov    %eax,0x4(%esp)
0x401429    cmpl   $0x6,0x4(%esp)
0x40142e    ja     0x4014a0
```

Just add

mov %eax,%ecx
mov %ecx,%eax

and break results

```
0x4015a0 ...
0x4015a5 call D
...      ...
```

```
0x401470 ...
0x401475 call F1
...      ...
```

```
0x4014f0 ...
0x4014f5 call F2
...      ...
```

```
0x4014a0 ...
0x4014a5 call F3
...      ...
```

```
0x401450 ...
0x401455 call F0
...      ...
```

```
0x4016d0 leave
0x4016d1 ret
```

## With IDA

- **Static (syntactic): too fragile**
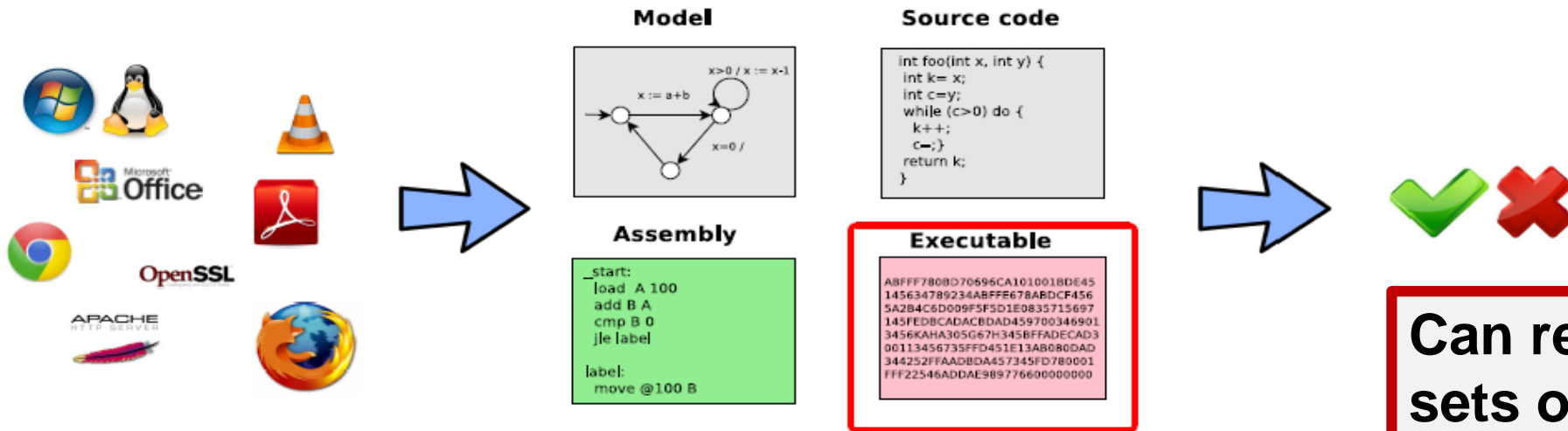- **Dynamic: too incomplete**

- **Malware deobfuscation is necessary**

- **Malware deobfuscation is highly challenging**

- **Standard tools are not enough – experts need better help!**

- **Static (syntactic): too fragile**
- **Dynamic: too incomplete**

# SOLUTION? BINARY-LEVEL SEMANTIC ANALYSIS

**Semantic preserved by obfuscation**

**Semantic tools** help make sense of binary
- Develop the next generation of binary-level tools !
- motto : leverage formal methods from safety critical systems

**Model**

**Source code**
```
int foo(int x, int y) {
    int k= x;
    int c=y;
    while (c>0) do {
        k++;
        c--;}
    return k;
}
```

**Assembly**
```
_start:
    load  A 100
    add B A
    cmp B 0
    jle label

label:
    move @100 B
```

**Executable**
```
A8FFF780BD70696CA101001BDE45
145634789234ABFFE678ABDCF456
5A2B4C6D009F5F5D1E0835715697
145FEDBCADACBDAD459700346901
3456KAHA305G67H345BFFADECAD3
00113456735FFD451E13AB080DAD
344252FFAADBDA457345FD780001
FFF22546ADDAE9B9776600000000
```

**Can reason about sets of executions**
- find rare events
- prove infeasibility

## Advantages
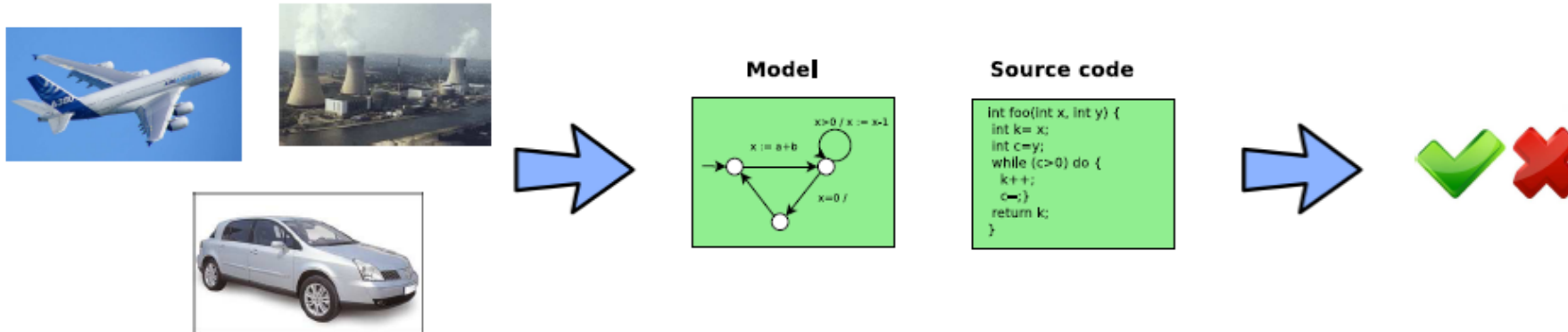- more robust than syntactic
- more thorough than dynamic

## Challenges
- source-level ↦ binary-level
- safety ↦ security
- many (complex) architectures

# *<En aparté>* ABOUT FORMAL METHODS

- Between Software Engineering and Theoretical Computer Science
- Goal = proves correctness in a mathematical way

**Success in safety-critical**



**Key concepts : $M \models \varphi$**

- $M$ : semantic of the program
- $\varphi$ : property to be checked
- $\models$ : algorithmic check

**Kind of properties**

- absence of runtime error
- pre/post-conditions
- temporal properties

# *<En aparté>* A DREAM COME TRUE … IN CERTAIN DOMAINS

**Industrial reality** in some **key areas**, especially safety-critical domains

- hardware, aeronautics [airbus], railroad [metro 14], smartcards, drivers [Windows], certified compilers [CompCert] and OS [Sel4], etc.
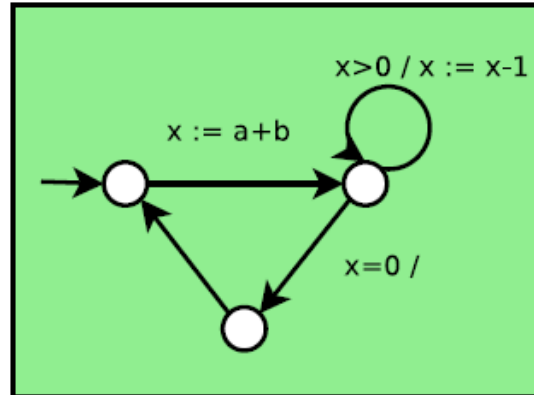
## Ex : Airbus

Verification of

- runtime errors [Astrée]

- functional correctness [Frama-C *]

- numerical precision [Fluctuat *]

- source-binary conformance [CompCert]

- ressource usage [Absint]

* : by CEA DILS/LSL

# NOW: BINARY-LEVEL ANALYSIS & OBFUSCATION

## Model



x := a+b

x>0 / x := x-1

x=0 /

## Source code

```
int foo(int x, int y) {
 int k= x;
 int c=y;
 while (c>0) do {
   k++;
   c--;}
 return k;
}
```

## Assembly

```
_start:
 load  A 100
 add B A
 cmp B 0
 jle label

label:
 move @100 B
```

## Executable

ABFFF780BD70696CA101001BDE45
145634789234ABFFE678ABDCF456
5A2B4C6D009F5F5D1E0835715697
145FEDBCADACBDAD459700346901
3456KAHA305G67H345BFFADECAD3
00113456735FFD451E13AB080DAD
344252FFAADBDA457345FD780001
FFF22546ADDAE989776600000000

## Low-level semantics of data

- machine arithmetic, bit-level operations, untyped memory
- ▶ difficult for any state-of-the-art formal technique

## Low-level semantics of control

- no distinction data / instructions, dynamic jumps (`jmp eax`)
- no (easy) syntactic recovery of Control-Flow Graph (CFG)
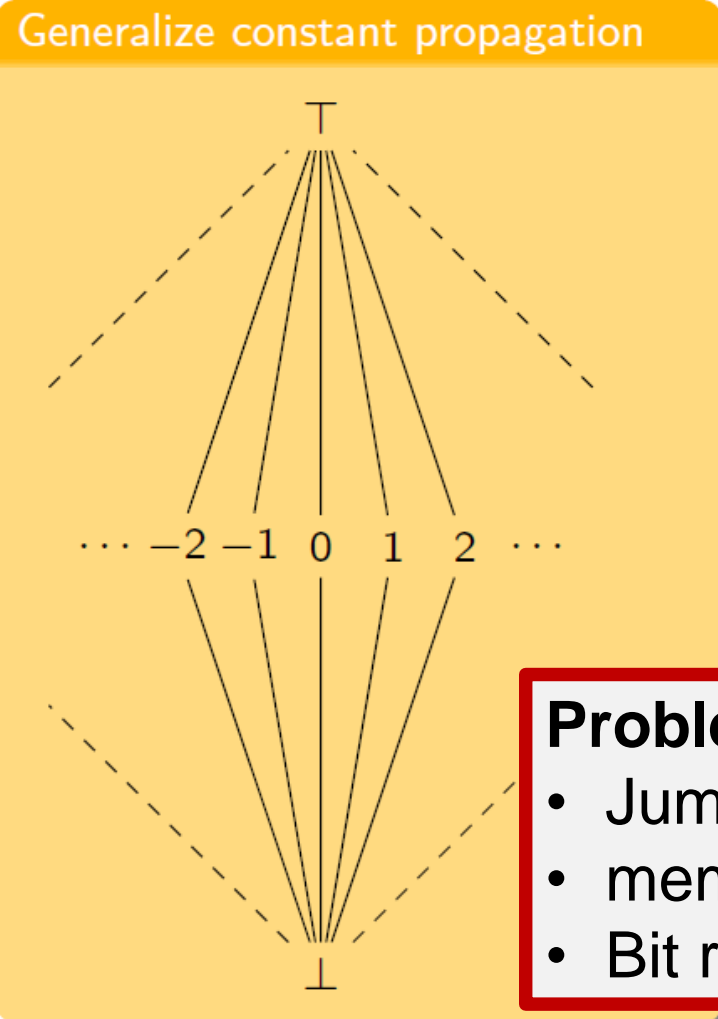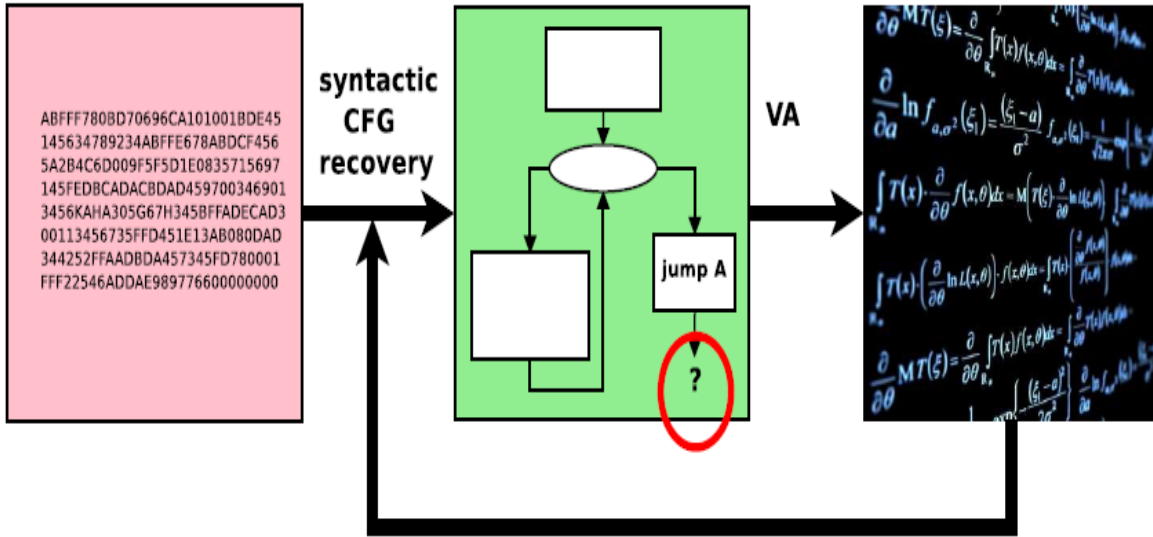- ▶ violate an implicit prerequisite for most formal techniques

## Diversity of architectures and instruction sets

- support for many instructions, modelling issues
- ▶ tedious, time consuming and error prone

**Wanted**
- robustness
- precision
- scale

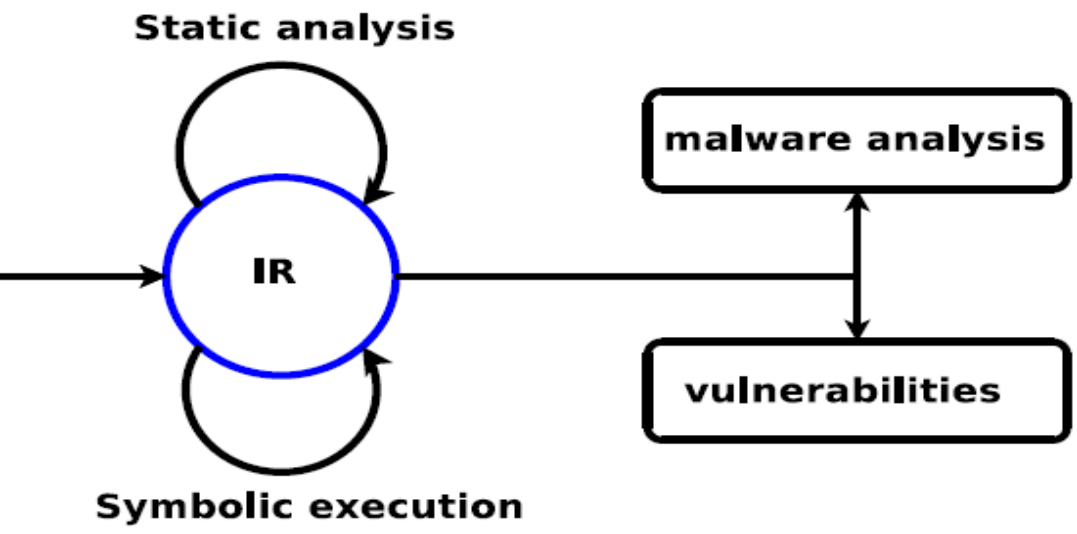# *<En aparté>* STATIC SEMANTIC ANALYSIS IS VERY VERY HARD ON BINARY CODE



syntactic CFG recovery

VA

jump A

?

## Framework : abstract interpretation

- notion of abstract domain
  $\bot, \top, \sqcup, \sqcap, \sqsubseteq, eval^{\#}$
- more or less precise domains
  . intervals, polyhedra, etc.
- fixpoint until stabilization

## Generalize constant propagation

$$\top$$

$$\cdots\ -2\ -1\ \ 0\ \ 1\ \ 2\ \cdots$$

$$\bot$$

## Problems
- Jump eax
- memory
- Bit resoning

# OUR APPROACH: BINSEC



**x86**

```
ABFFF780BD70696CA101001BDE45
145634789234ABFFE678ABDCF456
5A2B4C6D009F5F5D1E0835715697
145FEDBCADACBDAD459700346901
3456KAHA305G67H345BFFADECAD3
00113456735FFD451E13AB080DAD
344252FFAADBDA457345FD780001
FFF22546ADDAE989776600000000
```
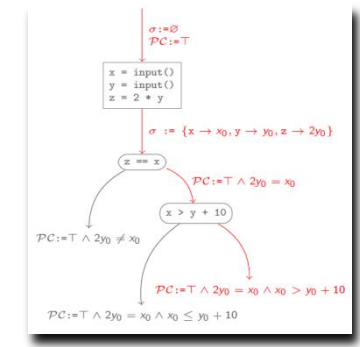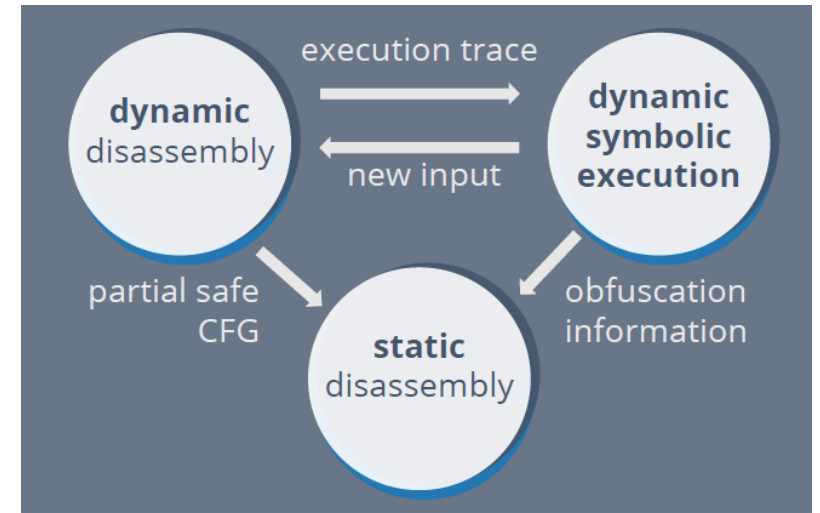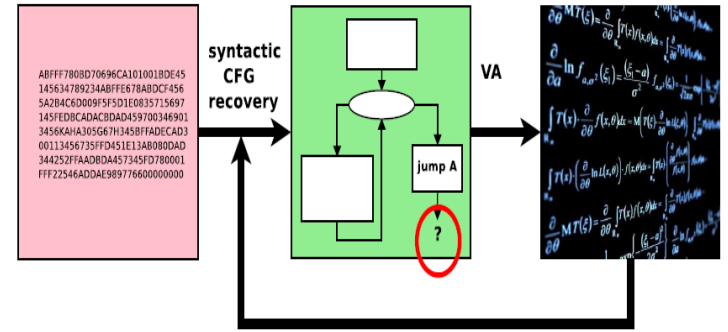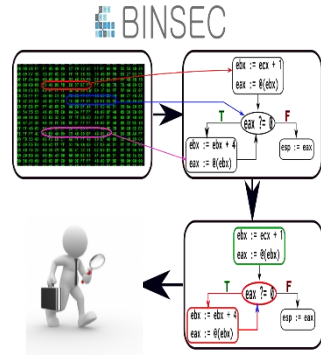
**ARM**

```
ABFFF780BD70696CA101001BDE45
145634789234ABFFE678ABDCF456
5A2B4C6D009F5F5D1E0835715697
145FEDBCADACBDAD459700346901
3456KAHA305G67H345BFFADECAD3
00113456735FFD451E13AB080DAD
344252FFAADBDA457345FD780001
FFF22546ADDAE989776600000000
```

**...**

```
ABFFF780BD70696CA101001BDE45
145634789234ABFFE678ABDCF456
5A2B4C6D009F5F5D1E0835715697
145FEDBCADACBDAD459700346901
3456KAHA305G67H345BFFADECAD3
00113456735FFD451E13AB080DAD
344252FFAADBDA457345FD780001
FFF22546ADDAE989776600000000
```
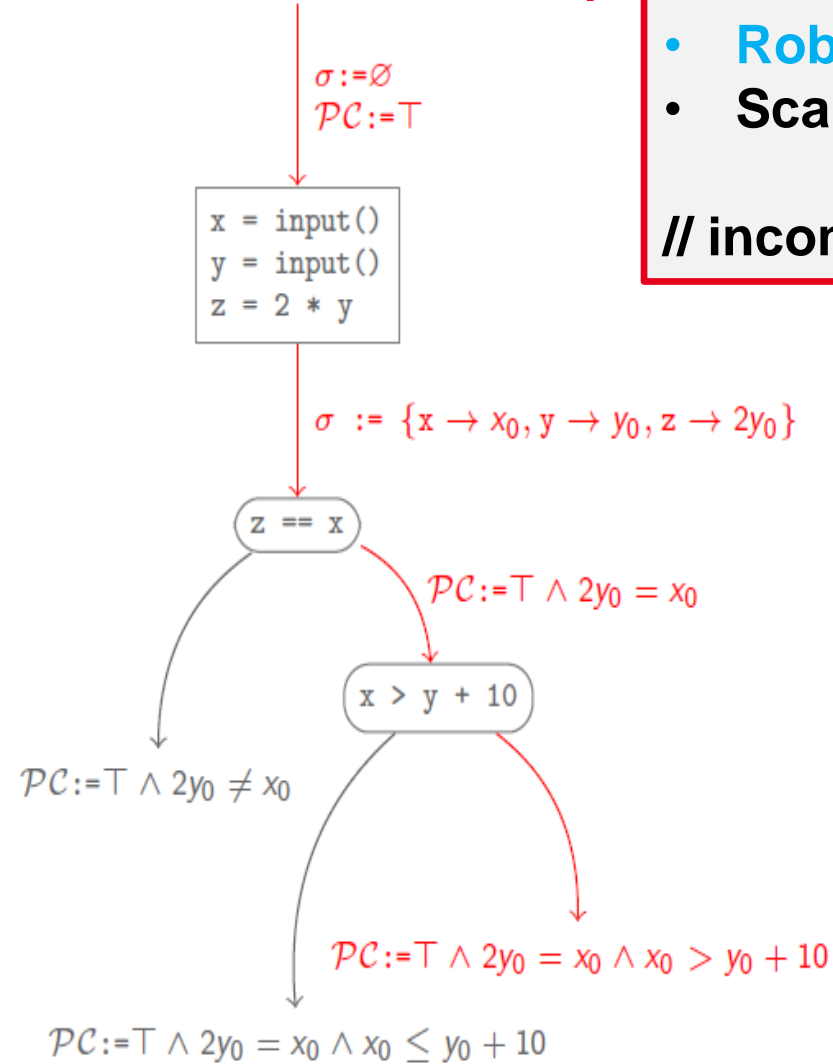
**Static analysis**

**IR**

**Symbolic execution**

**malware analysis**

**vulnerabilities**

- lhs := rhs
- goto addr, goto expr
- ite(cond)? goto addr :
- assume, assert, nondet

# KEY: DYNAMIC SYMBOLIC EXECUTION
## (DSE, Godefroid 2005)

```
int main () {
    int x = input();
    int y = input();
    int z = 2 * y;
    if (z == x) {
        if (x > y + 10)
            failure;
    }
    success;
}
```

- given a path of the program
- automatically find input that follows the path
- then, iterate over all paths

$\sigma := \varnothing$
$\mathcal{PC} := \top$

```
x = input()
y = input()
z = 2 * y
```

$\sigma := \{x \rightarrow x_0, y \rightarrow y_0, z \rightarrow 2y_0\}$

$z == x$

$\mathcal{PC} := \top \wedge 2y_0 = x_0$

$x > y + 10$

$\mathcal{PC} := \top \wedge 2y_0 \neq x_0$

$\mathcal{PC} := \top \wedge 2y_0 = x_0 \wedge x_0 > y_0 + 10$

$\mathcal{PC} := \top \wedge 2y_0 = x_0 \wedge x_0 \leq y_0 + 10$

# DYNAMIC SYMBOLIC EXECUTION CAN HELP (Debray, Kruegel, …)

```
int main () {
    int x = input();
    int y = input();
    int z = 2 * y;
    if (z == x) {
        if (x > y + 10)
            failure;
    }
    success;
}
```

- given a path of the program
- automatically find input that follows the path
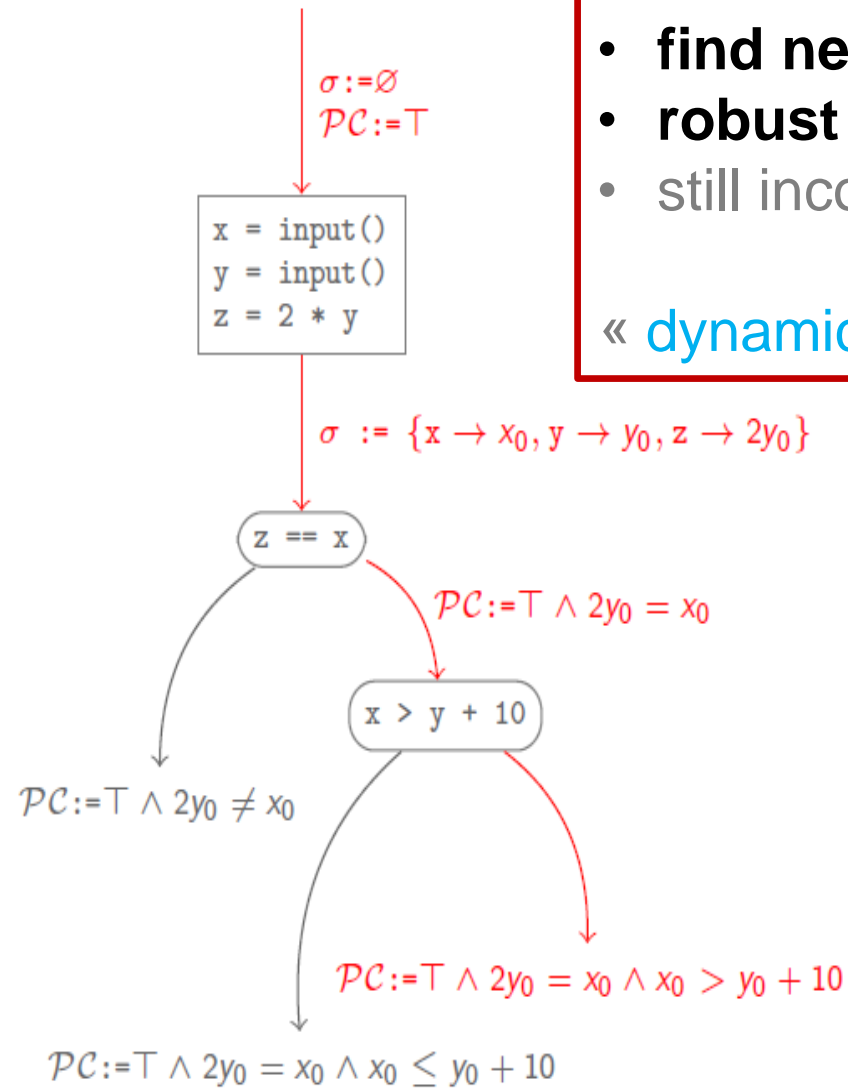- then, iterate over all paths

**For deobfuscation**
- **find new real paths**
- **robust**
- still incomplete

« dynamic analysis on steroids »

$\sigma := \varnothing$
$\mathcal{PC} := \top$

```
x = input()
y = input()
z = 2 * y
```

$\sigma := \{x \to x_0, y \to y_0, z \to 2y_0\}$

z == x

$\mathcal{PC} := \top \land 2y_0 = x_0$

$\mathcal{PC} := \top \land 2y_0 \neq x_0$

x > y + 10

$\mathcal{PC} := \top \land 2y_0 = x_0 \land x_0 > y_0 + 10$

$\mathcal{PC} := \top \land 2y_0 = x_0 \land x_0 \leq y_0 + 10$

```
                                    0x4013e0   push %ebp
                                    0x4014e1   mov %esp,%ebp
...                                 ...        ...
0x401419   mov    0xc(%esp),%eax    0x401430   mov    0x4(%esp),%eax
0x40141d   sub    $0x4,%eax         0x401434   shl    $0x2,%eax
0x401420   imul   0xc(%esp),%eax    0x401437   add    $0x40a064,%eax
0x401425   mov    %eax,0x4(%esp)    0x40143c   mov    (%eax),%eax
0x401429   cmpl   $0x6,0x4(%esp)    0x401441   mov    %eax,%ecx
0x40142e   ja     0x4014a0          0x401446   mov    %ecx,%eax
                                    0x40144b   jmp    *%eax
```

```
0x4015a0 ...        0x401470 ...        0x4014f0 ...        0x4014a0 ...        0x401450 ...
0x4015a5 call D     0x401475 call F1    0x4014f5 call F2    0x4014a5 call F3    0x401455 call F0
...      ...        ...      ...        ...      ...        ...      ...        ...      ...
```

```
0x4016d0 leave
0x4016d1 ret
```

With IDA + BINSEC

```
cmp eax ebx

cmc

jae ...
```

$$CF := (eax <_u ebx)$$
$$CF := \neg CF$$
$$if\ (\neg CF)\ goto\ ...$$

**Can recover useful semantic information**
- **More precise disassembly**
- **Exact semantic of instructions**
- **Input of interest**
- **…**

```
...                          0x40104d  mov al, input_buffer[ecx]
0x401045  call ReadFile      0x401053  xor al, 0x7d                    zf = 0     0x40107b  ...
0x40104b  xor ecx, ecx       0x401055  cmp al, data_str[ecx]                     ...
                             0x40105b  jnz 0x40107b                              You are failure

                     zf = 1                              ecx < 0x18

bunny_slope@flare-on.com     0x40105d  inc ecx
                             0x40105e  cmp ecx, 0x18      ecx = 0x18    0x401063  ...
                             0x401061  jl 0x40104d                      ...
                                                                       You are success
```
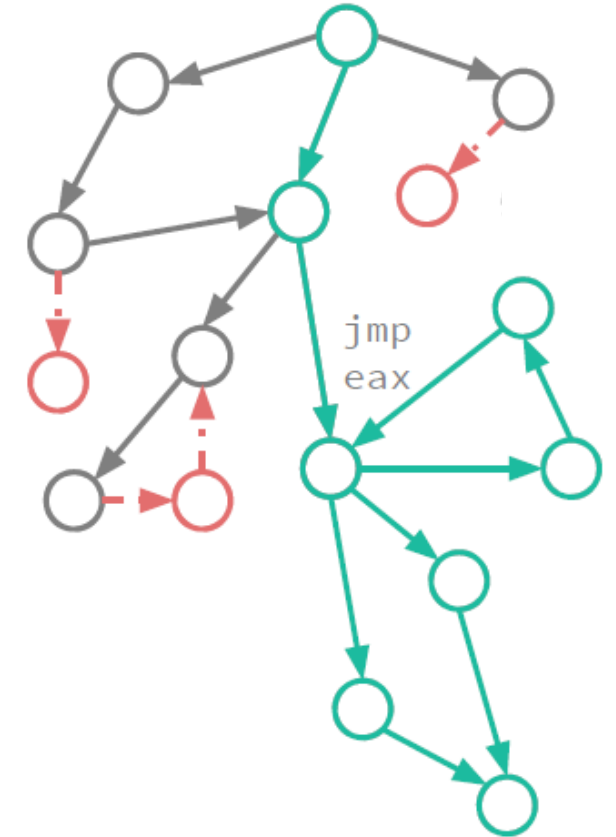
**Prove that something is always true (resp. false)**

**Many such issues in reverse**
- **is a branch dead?**
- **does the ret always return to the call?**
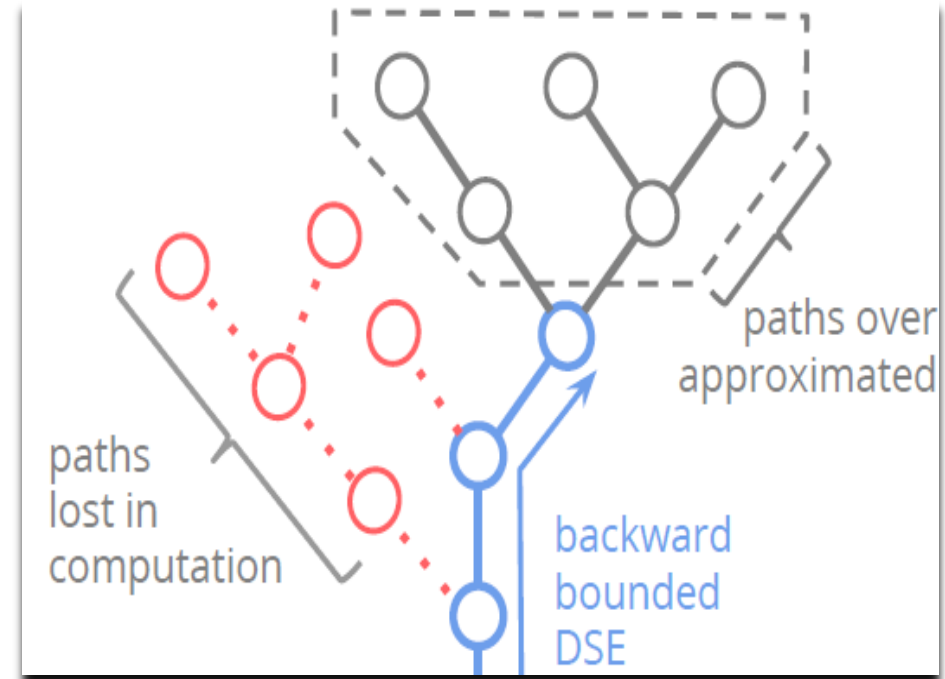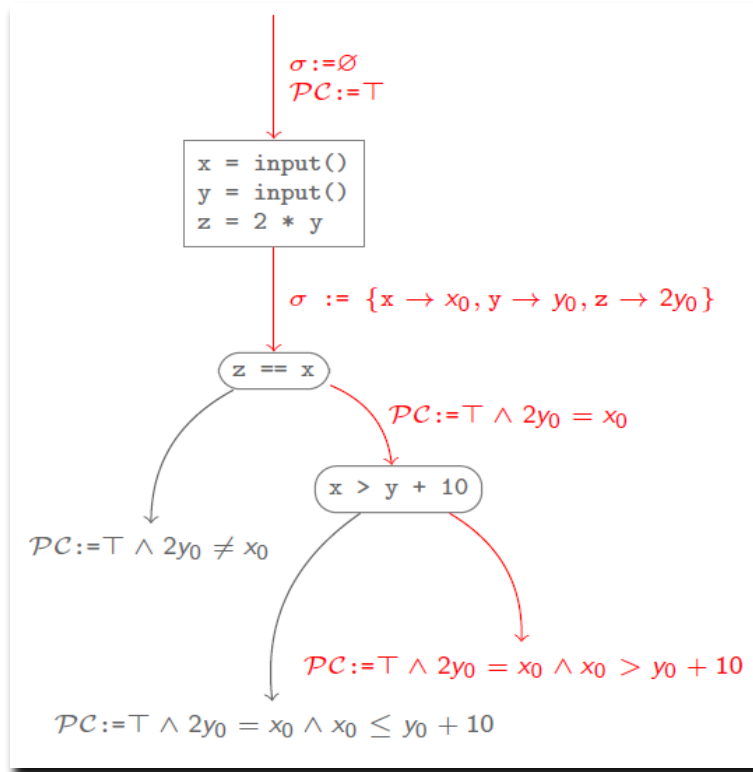- **have i found all targets of a dynamic jump?**

**And more**
- **does this malicious ret always go there?**
- **does this expression always evaluate to 15?**
- **does this self-modification always write this opcode?**
- **does this self-modification always rewrite this instr.?**
- **…**

jmp
eax

**Not addressed by DSE**
- **Cannot enumerate all paths**

# FORWARD & BACKWARD SYMBOLIC EXECUTION



| | (forward) DSE | bb-DSE |
|---|---|---|
| feasibility queries | 🟢 | 🔴 |
| infeasibility queries | 🔴 | 🟢 |
| scale | 🟠 | 🟢 |

- **Controlled experiments** (ground truth) ➡ **precision**

- **Large-scale experiment: packers** ➡ **scalability, robustness**

- **Case-study: X-tunnel malware** ➡ **usefulness**

- **Goal = assess the precision of the technique**
  - ground truth value

- **Experiment 1: opaque predicates (o-llvm)**
  - 100 core utils, 5x20 obfuscated codes
  - k=16: 3.46% error, no false negative
  - robust to k
  - efficient: 0.02s / query

- **Experiment 2: stack tampering (tigress)**
  - 5 obfuscated codes, 5 core utils
  - almost all genuine ret are proved (no false positive)
  - many malicious ret are proved « single-targets »

| k | OP (5556) ok | miss (FN) | Genuine (5183) ok | miss (FP) | TO | Error rate (FP+FN)/Tot (%) | Time (s) | avg/query (s) |
|---|---|---|---|---|---|---|---|---|
| 2 | 0 | 5556 | 5182 | 1 | 0 | 51.75 | 89 | 0.008 |
| 4 | 903 | 4653 | 5153 | 30 | 0 | 43.61 | 96 | 0.009 |
| | | | | | | | 120 | 0.011 |
| | | | | | | | 152 | 0.014 |
| | | | | | | | 197 | 0.018 |
| | | | | | | | 272 | 0.025 |
| | | | | | | | 384 | 0.036 |
| 32 | 5552 | 4 | 4579 | 604 | 25 | 5.66 | 699 | 0.065 |
| 40 | 5548 | 8 | 4523 | 660 | 39 | 6.22 | 1145 | 0.107 |
| 50 | 5544 | 12 | 4458 | 725 | 79 | 6.86 | 2025 | 0.189 |

**• Very precise résults**
**• Seems efficient**

| Sample | runtime genuine #ret [†] | proved genuine | proved a/d | runtime violation #ret [†] | proved a/d | proved single |
|---|---|---|---|---|---|---|
| *obfuscated programs* | | | | | | |
| simple-if | 6 | 6 | 6/0 | 9 | 0/0 | 8 |
| bin-search | 15 | 15 | 15/0 | 25 | 0/0 | 24 |
| bubble-sort | 6 | 6 | 6/0 | 15 | 0/1 | 13 |
| mat-mult | 31 | 31 | 31/0 | 69 | 0/0 | 68 |
| huffman | 19 | 19 | 19/0 | 23 | 0/3 | 19 |
| *non-obfuscated programs* | | | | | | |
| ls | 30 | 30 | 30/0 | 0 | - | - |
| dir | 35 | 35 | 35/0 | 0 | - | - |
| mktemp | 21 | 20 | 20/0 | 0 | - | - |
| od | 21 | 21 | 21/0 | 0 | - | - |
| vdir | 49 | 43 | 43/0 | 0 | - | - |

# CASE-STUDY: PACKERS



| packers | trace len. | #proc | #th | #SMC | opaque predicates OK | OP | call stack tampering OK | tamper |
|---|---|---|---|---|---|---|---|---|
| ACProtect v2.0 | 1.8M | 1 | 1 | | | 159 | 0 | 48 |
| ASPack v2.12 | 377K | 1 | | | | 24 | 11 | 6 |
| Crypter v1.12 | 1.1M | 1 | | | | 24 | 125 | 78 |
| Expressor | 635K | 1 | 1 | | | | 14 | 0 |
| FSG v2.0 | 68k | 1 | 1 | | | | 6 | 0 |
| Mew | 59K | 1 | 1 | 1 | 28 | 1 | 6 | 1 |
| PE Lock | 2.3M | 1 | 1 | 6 | 95 | 90 | 4 | 3 |
| RLPack | 941K | 1 | 1 | | | | 14 | 0 |
| TELock v0.51 | 406K | 1 | 1 | | | | 3 | 1 |
| Upack v0.39 | 711K | 1 | 1 | | | | 7 | 1 |

The technique scale on significant traces

Many true positives. Some packers are using it intensively

Packers using ret to perform the final tail transition to the entrypoint

**Packers: legitimate software protection tools**
**(basic malware: the sole protection)**

# CASE-STUDY: PACKERS (fun facts)

Several of the tricks detected by the analysis

Obsidium JD Pack PE Lock WinUpack Expressor PE Compact Armadillo Packman EP Protector ACProtect TELock SVk Yoda's Crypter Mew Neolite UPX MoleBox FSG Upack Crypter Yoda's Protector ASPack BoxedApp Petite nPack PE Spin Enigma Setisoft Themida RLPack Mystic VMProtect

**OP in ACProtect**

| | | |
|---|---|---|
| 1018f7a | js | 0x1018f92 |
| 1018f7c | jns | 0x1018f92 |

(and all possible variants
ja/jbe, jp/jnp, jo/jno..)

**OP in Armadillo**

| | | |
|---|---|---|
| 10330ae | xor | ecx, ecx |
| 10330b0 | jnz | 0x10330ca |

**CST in ACProtect**

| | |
|---|---|
| 1001000 | push 16793600 |
| 1001005 | push 16781323 |
| 100100a | ret |
| 100100b | ret |

**CST in ACProtect**

| | | |
|---|---|---|
| 1004328 | call | 0x1004318 |
| 1004318 | add | [esp], 9 |
| 100431c | ret | |

**CST in ASPack**

| | | |
|---|---|---|
| 10043a9 | mov | [ebp+0x3a8], eax |
| 10043af | popa | 0x10043bb at runtime |
| 10043b0 | jnz | 0x10043ba |
| | Enter SMC Layer 1 | |
| 10043ba | push | 0x10011d7 |
| 10043bf | ret | |

**OP (decoy) in ASPack**

```
10040fe: mov bl, 0x0
10041c0: cmp bl, 0x1          0x10040ff
10041c3: jnz 0x1004163        at runtime
```

ZF = 0        ZF = 1

```
1004163: jmp 0x100416d
[...]
```

```
1004105: inc [ebp+0xec]
[...]
```

**Nicknames:** APT28, Fancy Bear, Sofacy, Sednit, Pawn Storm

Ministry of **Defense** (France)

First seen

Government **Officials** (Poland)

**NATO**, EU institution
**Bundestag** (Germany)

**TV5 Monde** (France)

**DNC** Democratic National Committee (US)

Political Activists (Russia)

2008    2011-2014    2015    2016

Office (RCE) CVE-2015-2424
Windows (LPE) CVE-2015-1701
Java (x2) CVE-2015-[2590,4902]
Flash (x2) CVE-2015-[3043,7645]
Flash + Windows 10 sandbox escape win32k.sys

**Two heavily obfuscated samples**
- **Many opaque predicates**

**Goal: detect & remove protections**
- Identify 50% of code as spurious
- Fully automatic, < 3h

Mark → Extract →

|  | C637 Sample #1 | 99B4 Sample #2 |
|---|---|---|
| #total instruction | **505,008** | **434,143** |
| #alive | +279,483 | +241,177 |

# CASE-STUDY: THE XTUNNEL MALWARE (fun facts)

- **Protection seems to rely only on opaque predicates**

- **Only two families of opaque predicates**      $7y^2 - 1 \neq x^2$      $\dfrac{2}{x^2 + 1} \neq y^2 + 3$

- **Yet, quite sophisticated**
  - original OPs
  - interleaving between payload and OP computation
  - sharing among OP computations
  - possibly long dependencies chains (avg 8.7, upto 230)

# SECURITY ANALYSIS: COUNTER-MEASURES (and mitigations)

- **Long dependecy chains (evading the bound k)**
  - Not always requires the whole chain to conclude!
  - Can use a more flexible notion of bound (data-dependencies, formula size)

- **Hard-to-solve predicates (causing timeouts)**
  - A time-out is already a valuable information
  - Opportunity to find infeasible patterns (then matching), or signatures
  - Tradeoff between performance penalty vs protection focus
  - Note: must be input-dependent, otherwise removed by standard DSE optimizations

- **Anti-dynamic tricks (fool initial dynamic recovery)**
  - Can use the appropriate mitigations
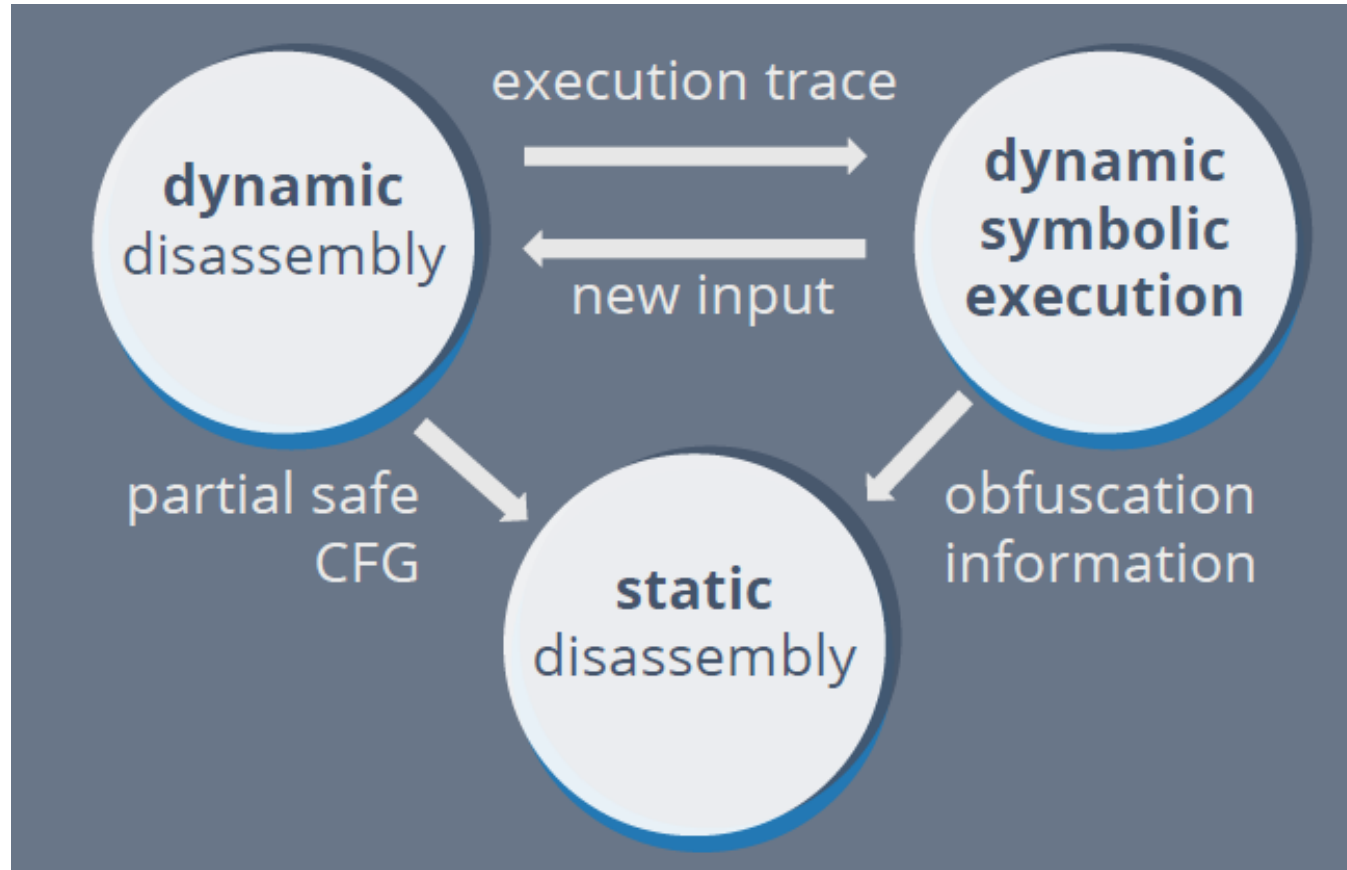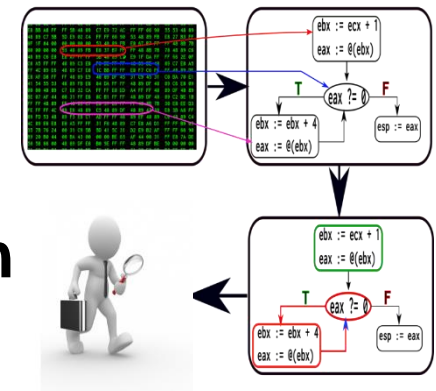  - Note: some tricks can be circumvent by symbolic reasoning

**Current state-of-the-art**
- **push the cat-and-mouse game further**
- **raise the bar for malware designers**

# SUMMARY

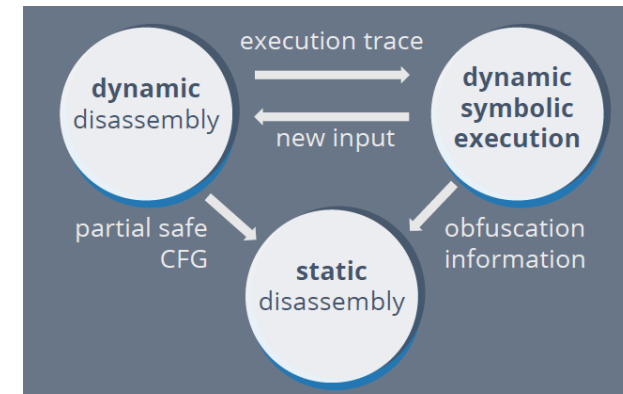| | Feasibility | Infeasibility | Efficient | Robust |
|---|---|---|---|---|
| **Static syntactic** | x | -- | OK | x |
| **Dynamic** | -- | x | OK | OK |
| | | | | |
| **DSE** | OK | x | x | OK |
| **BB-DSE** | x | OK | OK | OK |

# CONCLUSION & TAKE AWAY

- **A tour on the advantages of symbolic methods for deobfuscation**

- **Semantic analysis complements existing approaches**
  - Explore, prove infeasible, simplify
  - Open the way to fruitful combinations

- **Formal methods can be useful for malware, but must be adapted**
  - Need robustness and scalability!
  - Accept to lose both correctness & completeness – in a controlled way

- **Next Step**
  - Combines with user and learning!
  - Anti-anti-DSE