

BINSEC

BINary-level SECurity analaysis

Sébastien Bardin (coordinator)
Airbus group, CEA, IRISA, LORIA, Uni. Grenoble Alpes



Panorama de BINSEC

Cadre de travail

- Projet ANR INS, appel 2012
- Axes 1 (sécurité) et 2 (génie logiciel)
- Projet de recherche fondamentale sur 4 ans (2013-2017)

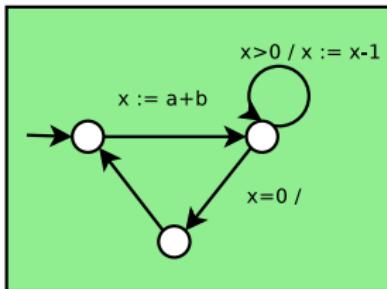
Sujet : Techniques formelles d'analyse de sécurité au niveau binaire

Partenaires : CEA (coordinateur), Airbus Group, INRIA Bretagne Atlantique, Université Grenoble Alpes - VERIMAG, Université de Lorraine - LORIA,

Participants : Sébastien Bardin, Frédéric Besson, Sandrine Blazy, Guillaume Bonfante, Richard Bonichon, Robin David, Adel Djoudi, Benjamin Farinier, Josselin Feist, Colas Le Guernic, Jean-Yves Marion, Laurent Mounier, Marie-Laure Potet, Than Dihn Ta, Franck Védrine, Pierre Wilke, Sara Zennou

Binary-level software analysis

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

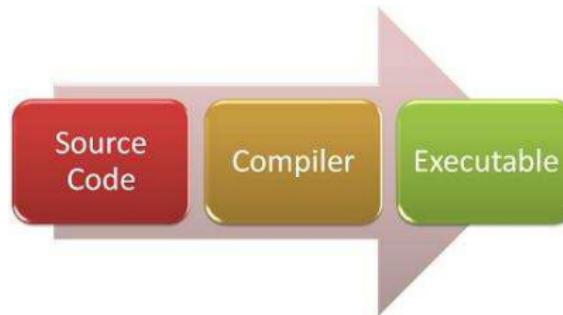
```
ABFFF780BD70696CA101001BDE45  
145634789234ABFFE678ABDCF456  
5A2B4C6D009F5F5D1E0835715697  
145FEDBCADACBDAD459700346901  
3456KAHA305G67H345BFFADECAD3  
00113456735FFD451E13AB080DAD  
344252FFAADBDA457345FD780001  
FFF22546ADDAE989776600000000
```

What for? (1)



How much do you trust your external components?

What for ? (2)



How much do you trust your compiler ?

What for? (2)

A perfectly legitimate compiler-introduced bug

```
void getPassword(void) {  
    char pwd [64];  
    if (GetPassword(pwd,sizeof(pwd))) {  
        /* checkpassword */  
    }  
    memset(pwd,0,sizeof(pwd));  
}
```

- Optimizing compilers may remove dead code
- `pwd` never accessed after `memset`
- **Thus can be safely removed**
- **And allows the password to stay longer in memory**

Mentioned in OpenSSH CVE-2016-0777

What for? (3)



A screenshot of a terminal window displaying binary data in hex format. The data consists of two columns of hex digits. The first column includes addresses (e.g., 00000000, 00000001, etc.) and the second column contains various hex values such as 33, 43, 94, b6, 00, cc, 96, dd, 2a, 32, 18, c7, a3, cf, 5a, 31, d9, 08, 48, db, 52, 62, e1, 0f, 30, 29, and 3. The background of the terminal is light blue.



Is it Stuxnet?

Binary-level security analysis

Several major security analyses are performed at byte-level

- vulnerability analysis [exploit finding]
- malware dissection and detection [deobfuscation]

State-of-the-technique

- very skilled experts, many efforts and **basic tools**
- dynamic analysis : gdb, fuzzing [easy to miss behaviours]
- static analysis : objdump, IDA Pro [easy to get fooled]

43	93	2a	52	c0	3
43	9c	32	62	08	2
43	10	32	48	4	
43	10	18	80	48	
d	94	9b	e1	db	0
b6	1d	07	a1	0f	5a
00	0e	a3	0f	5a	
84	cc	96	cf	31	
96	25	c9	dd	d9	
43	01	45	f0	29	



Binary-level security analysis

Several major security analyses are performed at byte-level

- vulnerability analysis [exploit finding]
- malware dissection and detection [deobfuscation]

State-of-the-technique

- very skilled experts, many efforts and **basic tools**
- dynamic analysis : gdb, fuzzing [easy to miss behaviours]
- static analysis : objdump, IDA Pro [easy to get fooled]

55	53	52	52	00	3
5c	2a	62	08	3	
32	62	48	48	48	
43	10	80	48	48	
94	9b	18	00	00	
d	94	e1	db	0	
b6	1d	07	0f	5a	
00	0e	a3	0f	5a	
84	cc	96	cf	31	
96	25	c9	dd	d9	
44	01	15	f0	29	



- goal : develop advanced tools for experts
- motto : leverage formal methods from safety critical systems

Not so easy!

Source-level \mapsto Binary-level

- semantics ?
- how to recover high-level structure ?

Safety \mapsto Security

- vulnerabilities : from warnings to exploitable bugs
- malware : from managed code to nasty code

7f	33	9c	2a	52	0d	3
0	7a	9c	2a	62	08	3
0	43	10	32	80	48	2
0	94	9b	18	80	4b	9
d	94	9b	18	80	4b	9
b6	1d	c7	e1	db	9	
00	0e	a3	0f	5a		
84	cc	96	cf	31		
96	25	c9	dd	d9		
dd	01	b5	f0	29		



Challenge : modelling

Instruction Prefixes	Opcode	ModR/M	SIB	Displacement	Immediate
Up to four prefixes of 1 byte each (optional)	1-, 2-, or 3-byte opcode	1 byte (if required)	1 byte (if required)	Address displacement of 1, 2, or 4 bytes or none	Immediate data of 1, 2, or 4 bytes or none

7	6 5	3 2	0	7	6 5	3 2	0
Mod	Reg/ Opcode	R/M		Scale	Index	Base	

Example of x86

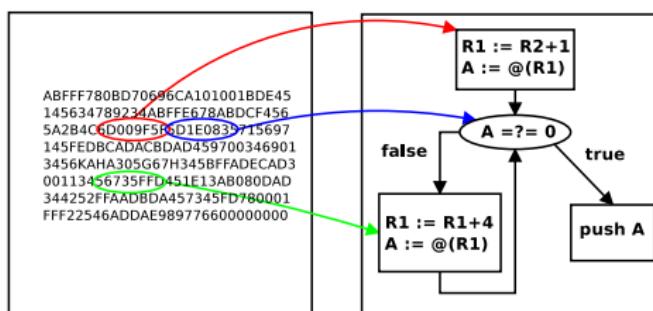
- more than 1,000 instructions
 - . ≈ 400 basic
 - . + float, interrupts, mmx
- many side-effects
- error-prone decoding
 - . addressing modes, prefixes, ...

Effective Address	Mod/R/M	Value of Mod	Byte
[EAX]	00	000	00
[ECX]	01	010	01
[EDX]	02	020	02
[EBX]	03	030	03
[ESI]	04	040	04
[EDI]	05	050	05
[disp32]	06	060	06
[TEST]	07	070	07
[ED1]	11	111	1F
[EAX]+disp8	00	0000	00
[ECX]+disp8	01	0100	01
[EDX]+disp8	02	0200	02
[EBX]+disp8	03	0300	03
[ESI]+disp8	04	0400	04
[EDI]+disp8	05	0500	05
[TEST]+disp8	06	0600	06
[ED1]+disp8	07	0700	07
[EAX]+disp32	00	00000000	00
[ECX]+disp32	01	01000000	01
[EDX]+disp32	02	02000000	02
[EBX]+disp32	03	03000000	03
[ESI]+disp32	04	04000000	04
[EDI]+disp32	05	05000000	05
[TEST]+disp32	06	06000000	06
[ED1]+disp32	07	07000000	07
AL/CL/EDX/ST1/MMX/XMM0	11	110	00
CL/ECX/EDX/ST1/MMX/XMM0	001	C0	D1
DL/DX/EDX/ST2/MMX/XMM0	010	C2	D2
BL/BX/EDX/ST3/MMX/XMM0	011	C3	D3
AH/CH/ESP/ST4/MMX/XMM0	100	C4	D4
CH/SH/ESP/ST5/MMX/XMM0	001	C5	D5
DH/SI/ESI/ST6/MMX/XMM0	110	C6	D6
BH/DI/EDI/ST7/MMX/XMM0	111	C7	D7

Challenge : safe CFG recovery

Input

- an executable code (array of bytes)
- an initial address
- a basic decoder : file \times address \mapsto instruction \times size



Output : (surapproximation of) the program CFG

[basic artifact for verif]

- problem : successors of $\langle \text{addr: goto a} \rangle$?

Challenges : vulnerabilities

Use-after-free (UaF) – CWE-416

- *dangling pointer* on **deallocated-then-reallocated** memory
- may lead to arbitrary data/code read, write or execution
- standard vulnerability in C/C++ applications (e.g. web browsers)
- challenges :
 - ▶ sequence of events, importance of aliasing
 - ▶ strongly depend on the implementation of `malloc` and `free`
 - ▶ scale, correctness (demonstrate bug), assess exploitability

```
1 char *login , *passwords;
2 login=(char *) malloc (...);
3 [...]
4 free(login); // login is now a dangling pointer
5 [...]
6 passwords=(char *) malloc (...); // ansi-C : fresh memory area
7 [...]
8 printf ("%s\n", login); // ansi-C : error, yet no particular security threat
```

Challenges : vulnerabilities

Use-after-free (UaF) – CWE-416

- *dangling pointer* on deallocated-then-reallocated memory
- may lead to arbitrary data/code read, write or execution
- standard vulnerability in C/C++ applications (e.g. web browsers)
- challenges
 - ▶ sequence of events, importance of aliasing
 - ▶ strongly depend on the implementation of malloc and free
 - ▶ scale, correctness (demonstrate bug), assess exploitability

```
char *login , *passwords;
2 login=(char *) malloc (...);
[...]
4 free(login); // login is now a dangling pointer
[...]
6 passwords=(char *) malloc (...); // may re-allocate memory of *login
[...]
8 printf("%s\n", login); // security threat : may print the passwords!
```

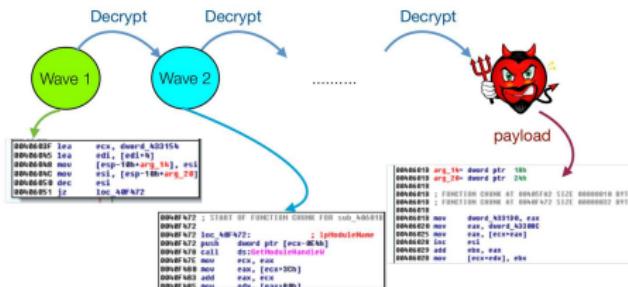
Challenge : malware

From managed code to nasty code

- recognize malware despite polymorphism
 - ▶ requires more robust notions of signatures

- recognize malware despite obfuscation
 - ▶ self-modifying code, opaque predicates, stack tampering, etc.

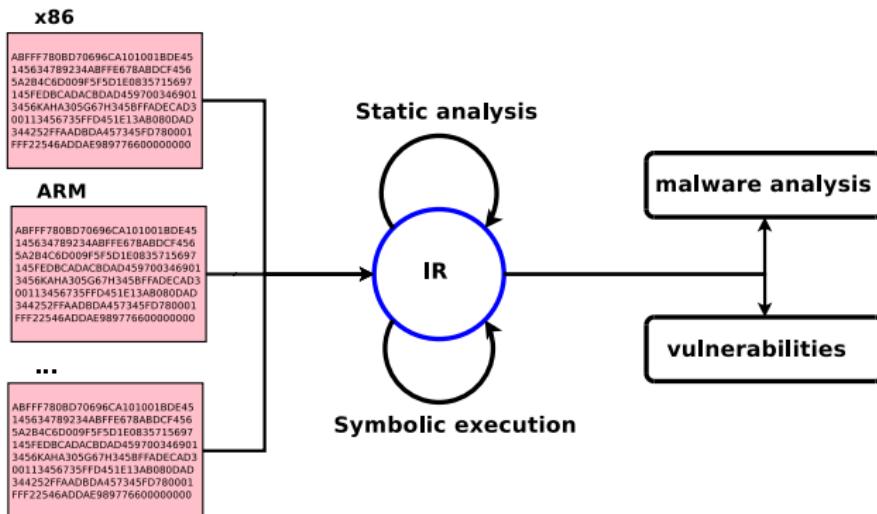
Context : x86-malware



A common protection scheme for malware
a SillyFDC run

Self-modifying program schema

BINSEC approach

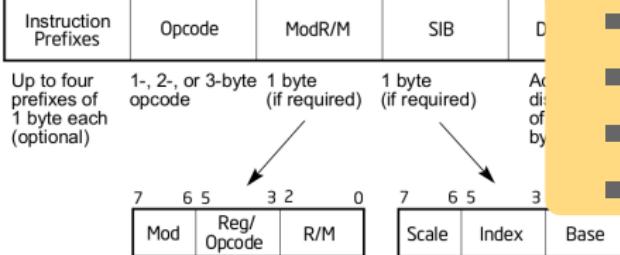


- leverage powerful methods from formal software analysis
- *pragmatic* formal methods (combination, tradeoffs, lose guarantees, etc.)
- common basic problems (models, analysis) + dedicated methods (vuln., malware)

Work and results

- Models and generic analysis [leader : irisa]
 - ▶ concise intermediate representation
 - ▶ low-level memory model for low-level programs [aplas14,itp15]
 - ▶ high-level structure recovery (cfg, conditions) [subm. fm16]
 - ▶ semantic exploration via symbolic execution [issta16]
- Dedicated security analysis [leaders : verimag, loria]
 - ▶ vulnerability : detection of use-after-free [j. virology 14,sub. woot16]
 - ▶ malware : semantic signatures [malware13-15,fps14]
 - ▶ deobfuscation : combined symbolic methods [ccs15, sub. ccs16]
- Open-source platform [leader : cea] [tacas15,saner16]
 - ▶ IR, frontend (loader, decoder, syntactic disassembly)
 - ▶ static analysis, symbolic execution
- Industrial evaluation [leader : airbus group]
 - ▶ industrial case-studies (iprint client, proftpd) and challenges (coreutils, vxheavens, samate, etc.)
 - ▶ Airbus Group : internal evaluation of GUEB
 - ▶ LORIA : startup around Gorille

Focus : modelling



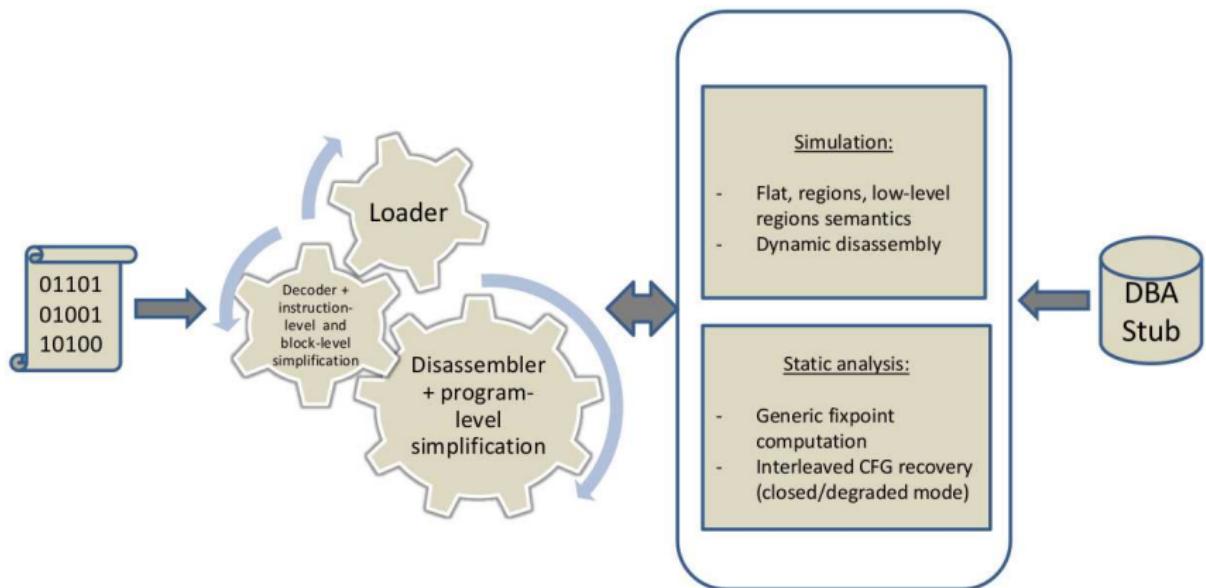
- `lhs := rhs`
- `goto addr, goto expr`
- `ite(cond)? goto addr : goto addr'`
- `assume, assert, nondet, malloc, free`

Intermediate Representation taken from [cav11]

- architecture independent
- really reduced set of instructions
 - . 9 instructions, less than 30 operators
- simple
 - clear semantic, no side-effect
- yet : no threads, no self-modif, no float

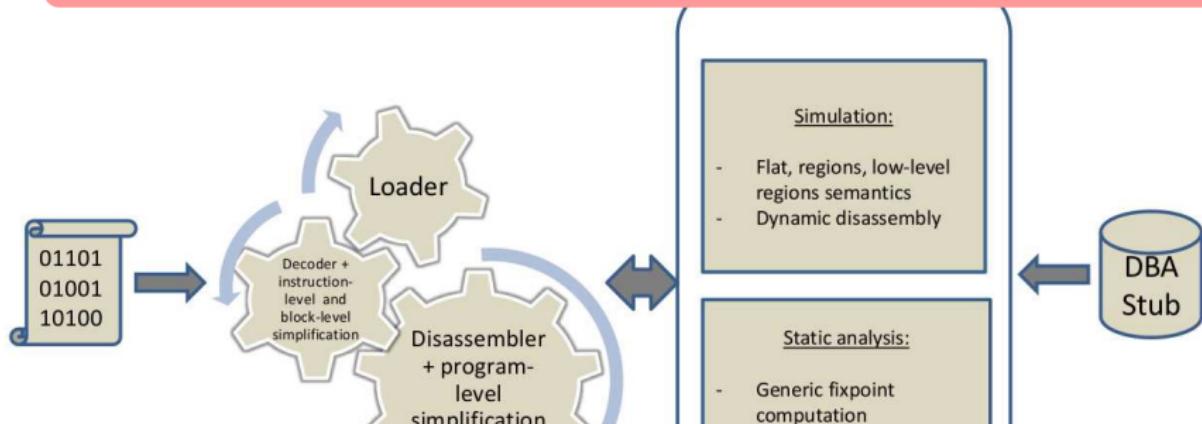
digit (opcode)	Mod/R/M	Value of Mod/R/M Byte (in Hex)
00	00	00
01	01	09
02	02	0A
03	03	0B
04	04	0C
05	05	0D
06	06	0E
07	07	0F
08	08	10
09	09	11
0A	0A	12
0B	0B	13
0C	0C	14
0D	0D	15
0E	0E	16
0F	0F	17
10	10	18
11	11	19
12	12	1A
13	13	1B
14	14	1C
15	15	1D
16	16	1E
17	17	1F
18	18	20
19	19	21
1A	1A	22
1B	1B	23
1C	1C	24
1D	1D	25
1E	1E	26
1F	1F	27
20	20	28
21	21	29
22	22	2A
23	23	2B
24	24	2C
25	25	2D
26	26	2E
27	27	2F
28	28	30
29	29	31
2A	2A	32
2B	2B	33
2C	2C	34
2D	2D	35
2E	2E	36
2F	2F	37
30	30	38
31	31	39
32	32	3A
33	33	3B
34	34	3C
35	35	3D
36	36	3E
37	37	3F
38	38	70
39	39	71
3A	3A	72
3B	3B	73
3C	3C	74
3D	3D	75
3E	3E	76
3F	3F	77
40	40	60
41	41	61
42	42	62
43	43	63
44	44	64
45	45	65
46	46	66
47	47	67
48	48	68
49	49	69
4A	4A	6A
4B	4B	6B
4C	4C	6C
4D	4D	6D
4E	4E	6E
4F	4F	6F
50	50	70
51	51	71
52	52	72
53	53	73
54	54	74
55	55	75
56	56	76
57	57	77
58	58	78
59	59	79
5A	5A	7A
5B	5B	7B
5C	5C	7C
5D	5D	7D
5E	5E	7E
5F	5F	7F
60	60	80
61	61	81
62	62	82
63	63	83
64	64	84
65	65	85
66	66	86
67	67	87
68	68	88
69	69	89
6A	6A	8A
6B	6B	8B
6C	6C	8C
6D	6D	8D
6E	6E	8E
6F	6F	8F
70	70	80
71	71	81
72	72	82
73	73	83
74	74	84
75	75	85
76	76	86
77	77	87
78	78	88
79	79	89
7A	7A	8A
7B	7B	8B
7C	7C	8C
7D	7D	8D
7E	7E	8E
7F	7F	8F
80	80	90
81	81	91
82	82	92
83	83	93
84	84	94
85	85	95
86	86	96
87	87	97
88	88	98
89	89	99
8A	8A	9A
8B	8B	9B
8C	8C	9C
8D	8D	9D
8E	8E	9E
8F	8F	9F
90	90	A0
91	91	A1
92	92	A2
93	93	A3
94	94	A4
95	95	A5
96	96	A6
97	97	A7
98	98	A8
99	99	A9
9A	9A	A0
9B	9B	A1
9C	9C	A2
9D	9D	A3
9E	9E	A4
9F	9F	A5
90	90	B0
91	91	B1
92	92	B2
93	93	B3
94	94	B4
95	95	B5
96	96	B6
97	97	B7
98	98	B8
99	99	B9
9A	9A	B0
9B	9B	B1
9C	9C	B2
9D	9D	B3
9E	9E	B4
9F	9F	B5
90	90	C0
91	91	C1
92	92	C2
93	93	C3
94	94	C4
95	95	C5
96	96	C6
97	97	C7
98	98	C8
99	99	C9
9A	9A	C0
9B	9B	C1
9C	9C	C2
9D	9D	C3
9E	9E	C4
9F	9F	C5
90	90	D0
91	91	D1
92	92	D2
93	93	D3
94	94	D4
95	95	D5
96	96	D6
97	97	D7
98	98	D8
99	99	D9
9A	9A	D0
9B	9B	D1
9C	9C	D2
9D	9D	D3
9E	9E	D4
9F	9F	D5
90	90	E0
91	91	E1
92	92	E2
93	93	E3
94	94	E4
95	95	E5
96	96	E6
97	97	E7
98	98	E8
99	99	E9
9A	9A	E0
9B	9B	E1
9C	9C	E2
9D	9D	E3
9E	9E	E4
9F	9F	E5
90	90	F0
91	91	F1
92	92	F2
93	93	F3
94	94	F4
95	95	F5
96	96	F6
97	97	F7
98	98	F8
99	99	F9
9A	9A	F0
9B	9B	F1
9C	9C	F2
9D	9D	F3
9E	9E	F4
9F	9F	F5

BINSEC platform



BINSEC platform

- loader ELF/PE, decoder x86
- 460/500 instructions : 380/380 "basic", 80/120 SIMD, no float/system
- prefixes : op size, addr size, repetition



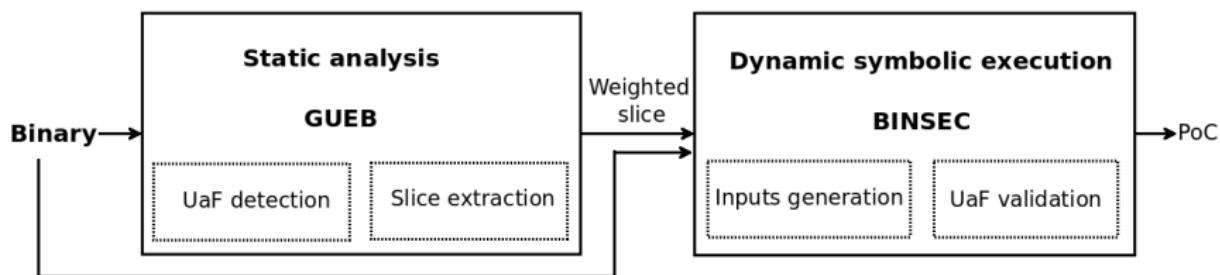
Basic services to build analysis on :

- Simulation
- Static analysis [safe CFG recovery] [tacas15, sub. fm16]
- Symbolic execution [flexible, optimized] [saner16, issta16, sub. fm16]

Finding *use-after-free* vulnerabilities

A *pragmatic two-step approach* implemented within the BINSEC platform :

- not complete, but scalable and correct in some cases



- **GUEB** : *scalable* lightweight static analysis (not sound, not complete)
→ produces a set of CFGs slices containing **potential** UaF
- **BINSEC/SE** : guided symbolic execution
→ *confirm* the UaF by finding **concrete** program inputs

Results

Several new vulnerabilities found

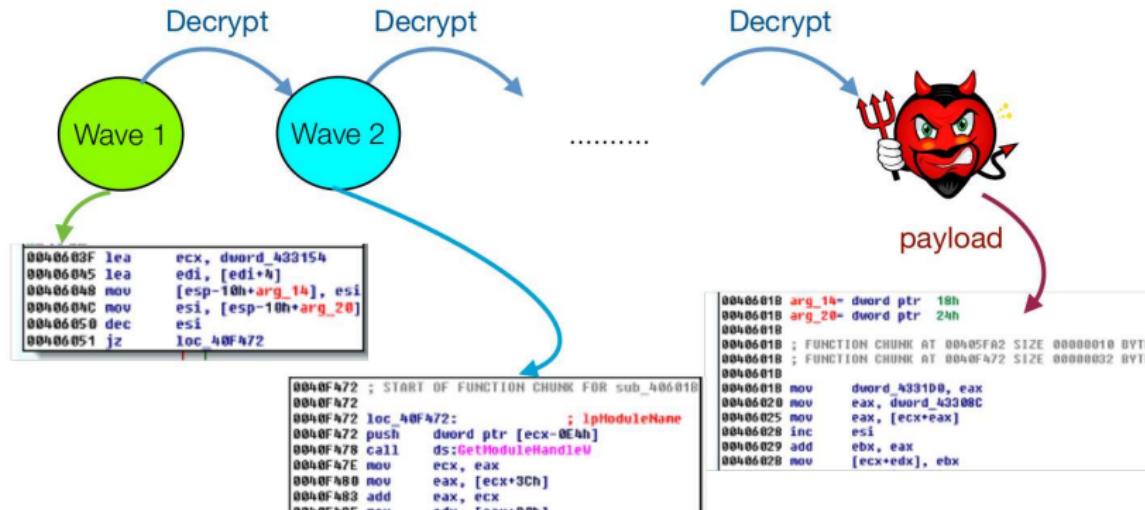
- GUEB + manual analysis [j. comp. virology 14]
 - ▶ tiff2pdf : CVE-2013-4232
 - ▶ openjpeg : CVE-2015-8871
 - ▶ gifcolor : CVE-2016-3177
 - ▶ accel-ppp
- GUEB + BINSE/SE [sefm16,sub. woot16]
 - ▶ Jasper JPEG-2000 : CVE-2015-5221

In progress

- third step : exploitability analysis of UaF
 - SMT-axiomatization of memory allocation strategies
- proFTPD industrial case study

Malware detection and Deobfuscation

Context : x86-malware



A common protection scheme for malware
a SillyFDC run

Self-modifying program schema

Malware detection and Deobfuscation

Context : x86-malware

Decrypt Decrypt Decrypt

Recognize malware despite polymorphism

- morphological analysis [malware13,fps14,malware15]
- experiments on Stuxnet, Flame, Duqu, Gauss, Waledec, Regin
- startup Simorfo since 2016

```

0040604C mov    esi, [esp-10h+arg_20]
00406050 dec    esi
00406051 jz     loc_404F72

```

```

00406018 arg_18= dword ptr 18h
00406018 arg_20= dword ptr 24h
00406018
00406018 ; FUNCTION CHUNK AT 00405FA2 SIZE 00000010 BYT
00406018 ; FUNCTION CHUNK AT 0040F472 SIZE 00000032 BYT
00406018

```

- standard anti-virus signatures \approx byte sequences
 - . easy to fool with slight rewriting (polymorphism)
- here : signatures = normalized subgraphs
 - . much harder to fool !
 - . efficient to compute / check
 - . in practice : add white lists

Malware detection and Deobfuscation

Context : x86-malware

Decrypt Decrypt Decrypt

Recognize malware despite obfuscation

- Codisasm [ccs15] : self-modification, code overlapping, anti-debug
- BINSEC/SE [saner16, sub. ccs16] : malware exploration, opaque predicate, stack tampering

00406048 mov [esp-10h+arg_18], esi
0040604C mov esi, [esp-10h+arg_20]
00406050 dec esi
00406051 jz loc_40F472

00406018 arg_18= dword ptr 18h
00406018 arg_20= dword ptr 24h
00406018 ; FUNCTION CHUNK AT 00405FA2 SIZE 00000010 BYT
00406018 ; FUNCTION CHUNK AT 0040F472 SIZE 00000032 BYT
00406018

- static analysis : not safe, complete, not robust to obfuscation
- dynamic analysis : safe, not complete, robust to obfuscation
- symbolic execution : best of both world
- combinations dynamic, static, symbolic

a OillyγDU TUI

Self-modifying program schema

Conclusion

Binary-level security analysis

- many applications
- many challenges

Our approach

- leverage powerful methods from formal software analysis
- *pragmatic* formal methods (combination, tradeoffs, etc.)
- common basic solutions (models, analysis) + application-directed methods

A few results

- vulnerability analysis [new CVE discovered]
- malware detection [startup]
- deobfuscation [applications on obfuscated malware]
- low-level memory regions [toward refining CompCert memory model]
- open-source platform [release soon]