

BINARY-LEVEL SECURITY: SEMANTIC ANALYSIS TO THE RESCUE

Sébastien Bardin (CEA LIST)

Joint work with

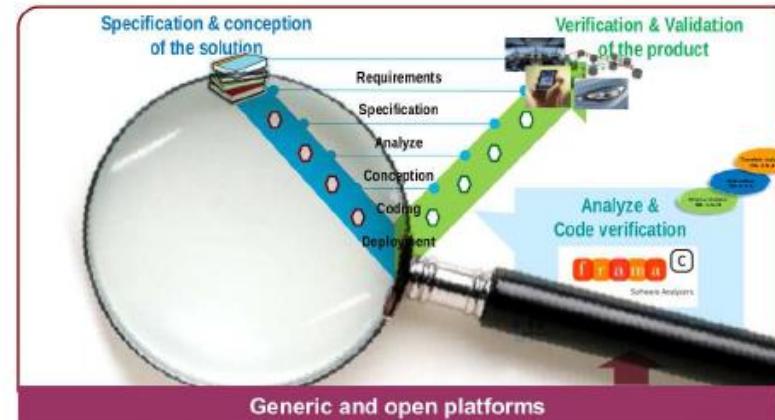
Richard Bonichon, Robin David, Adel Djoudi & many other people



ABOUT MY LAB @CEA

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



- **Binary-level security analysis: many applications, many challenges**
- **Standard techniques (dynamic, syntactic) not enough**
- **Formal methods can help ... but must be strongly adapted**
 - [Complement existing methods]
 - Need robustness, precision and scalability!
 - Acceptable to lose both correctness & completeness – in a controlled way
 - **New challenges and variations, many things to do!**
- **A tour on how formal methods can help**
 - *Explore and discover*
 - *Prove infeasibility or validity*
 - *Simplify* (not covered here)
 - *with Josselin Feist*
 - *with Robin David*
 - *with Jonathan Salwan*



- Why binary-level analysis?
- Some background on source-level formal methods
- The hard journey from source to binary
- A few case-studies
- Conclusion

- Focus mostly on Symbolic Execution
- Give hints for abstract Interpretation

Cover both

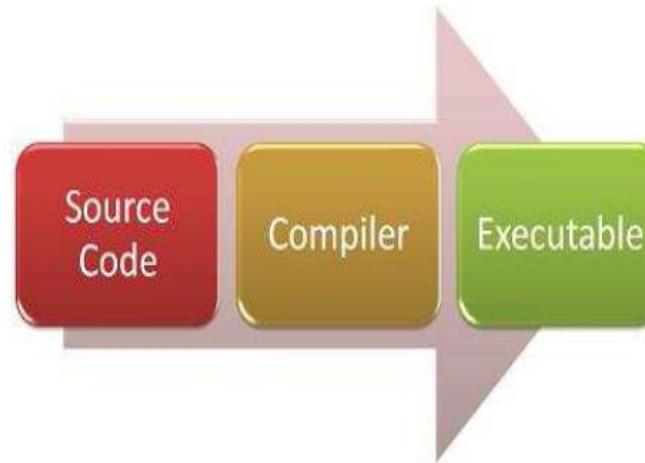
- vulnerability detection
- deobfuscation

- **Why binary-level analysis?**
- Some background on source-level formal methods
- The hard journey from source to binary
- A few case-studies
- Conclusion

BENEFITS



COTS



Is it Stuxnet ?

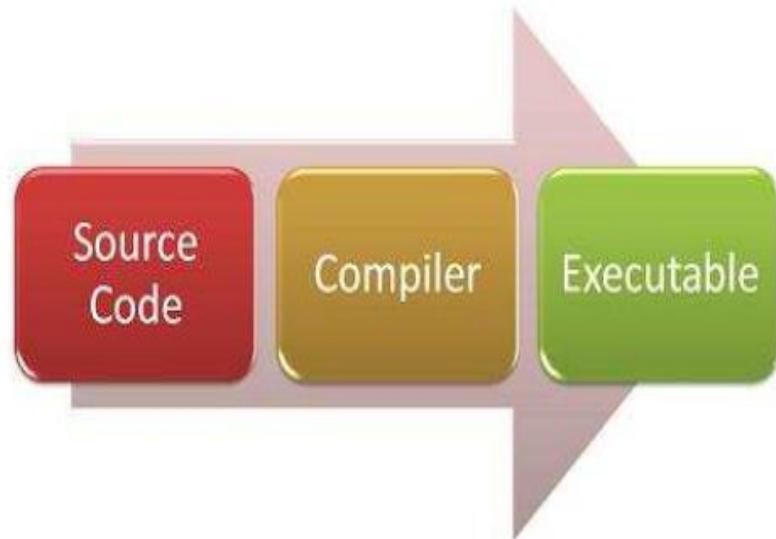
No source code

More precise analysis

Malware

What for: vulnerabilities, reverse (malware, legacy), protection evaluation, etc.

EXAMPLE: COMPILER BUG



- 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

Security bug introduced by a non-buggy compiler

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

OpenSSH CVE-2016-0777

Our goal here:
• Check the code after compilation

EXAMPLE: 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



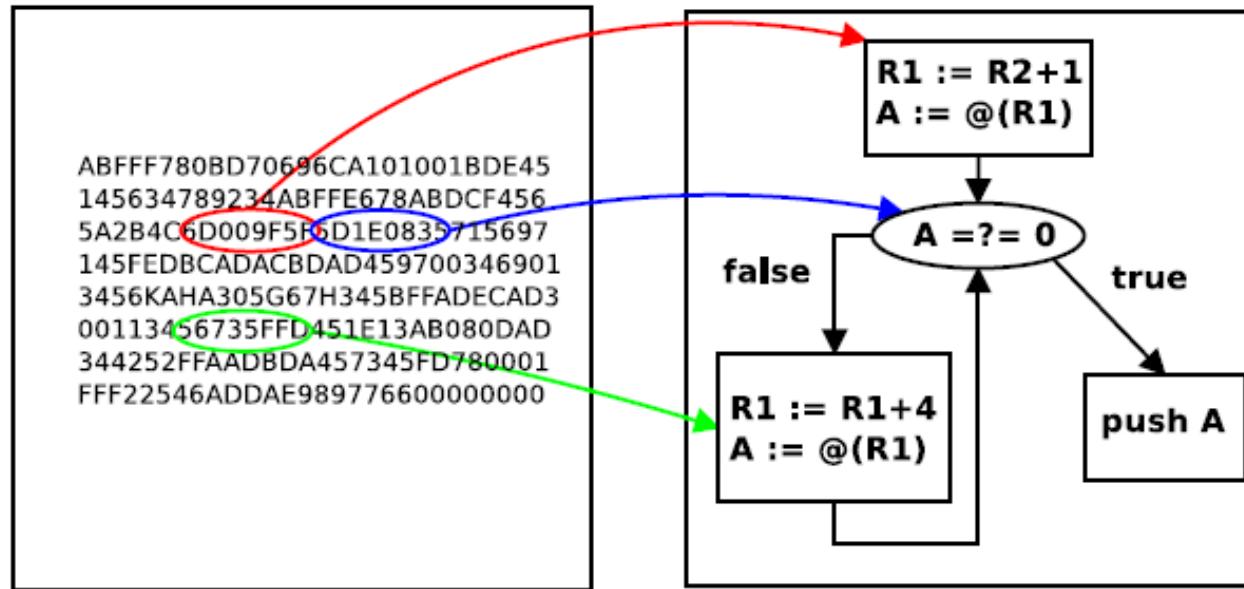
The day after: **malware comprehension**

- understand what has been going on
- mitigate, fix and clean
- improve defense



Highly challenging [obfuscation]

CHALLENGE: CORRECT DISASSEMBLY



Basic reverse problem

- aka model recovery
- aka CFG recovery

CAN BE TRICKY!

- code – data
- dynamic jumps (jmp eax)

Sections

.text

```
8D 4C 24 04 83 E4 F0 FF 71 FC 55 89 E5 53 51 83
EC 10 89 CB 83 EC 0C 6A 0A E8 A7 FE FF FF 83 C4
10 89 45 F0 8B 43 04 83 C0 04 8B 00 83 EC 0C 50
E8 C0 FE FF FF 83 C4 10 89 45 F4 83 7D F4 04 77
3B 8B 45 F4 C1 E0 02 05 98 85 04 08 8B 00 FF E0
C7 45 F4 00 00 00 00 EB 23 C7 45 F4 01 00 00 00
EB 1A C7 45 F4 02 00 00 00 EB 11 C7 45 F4 03 00
00 00 EB 08 C7 45 F4 04 00 00 00 90 83 EC 08 FF
75 F4 68 90 85 04 08 E8 29 FE FF FF 83 C4 10 8B
45 F4 8D 65 F8 59 5B 5D 8D 61 FC C3 66 90 66 90
66 90 66 90 90 55 57 31 FF 56 53 E8 85 FE FF FF
81 C3 89 12 00 00 83 EC 1C 8B 6C 24 30 8D B3 0C
FF FF FF E8 B1 FD FF FF 8D 83 08 FF FF FF 29 C6
C1 FE 02 85 F6 74 27 8D B6 00 00 00 00 8B 44 24
38 89 2C 24 89 44 24 08 8B 44 24 34 89 44 24 04
FF 94 BB 08 FF FF 83 C7 01 39 F7 75 DF 83 C4
1C 5B 5E 5F 5D C3 EB 0D 90 90 90 90 90 90 90 90 90
90 90 90 90 90 F3 C3 FF FF 53 83 EC 08 E8 13 FE
FF FF 81 C3 17 12 00 00 83 C4 08 5B C3 03 00 00
00 01 00 02 00 76 61 6C 3A 25 64 0A 00 AB 84 04
08 B4 84 04 08 BD 84 04 08 C6 84 04 08 CF 84 04
08 01 1B 03 3B 28 00 00 00 04 00 00 00 54 FD FF
```

.fini
.rodata

.eh_frame_hdr

■ code ■ dead bytes ■ global csts ■ strings ■ pointers ■ other

Code (Functions)

main

unknown

_libc_csu_init

unknown

_libc_csu_fini

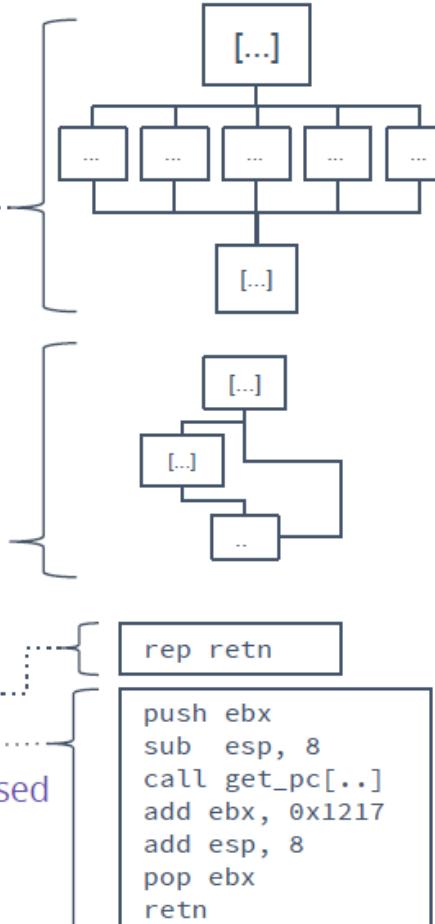
_term_proc

_fp_hw, _IO_stdin_used

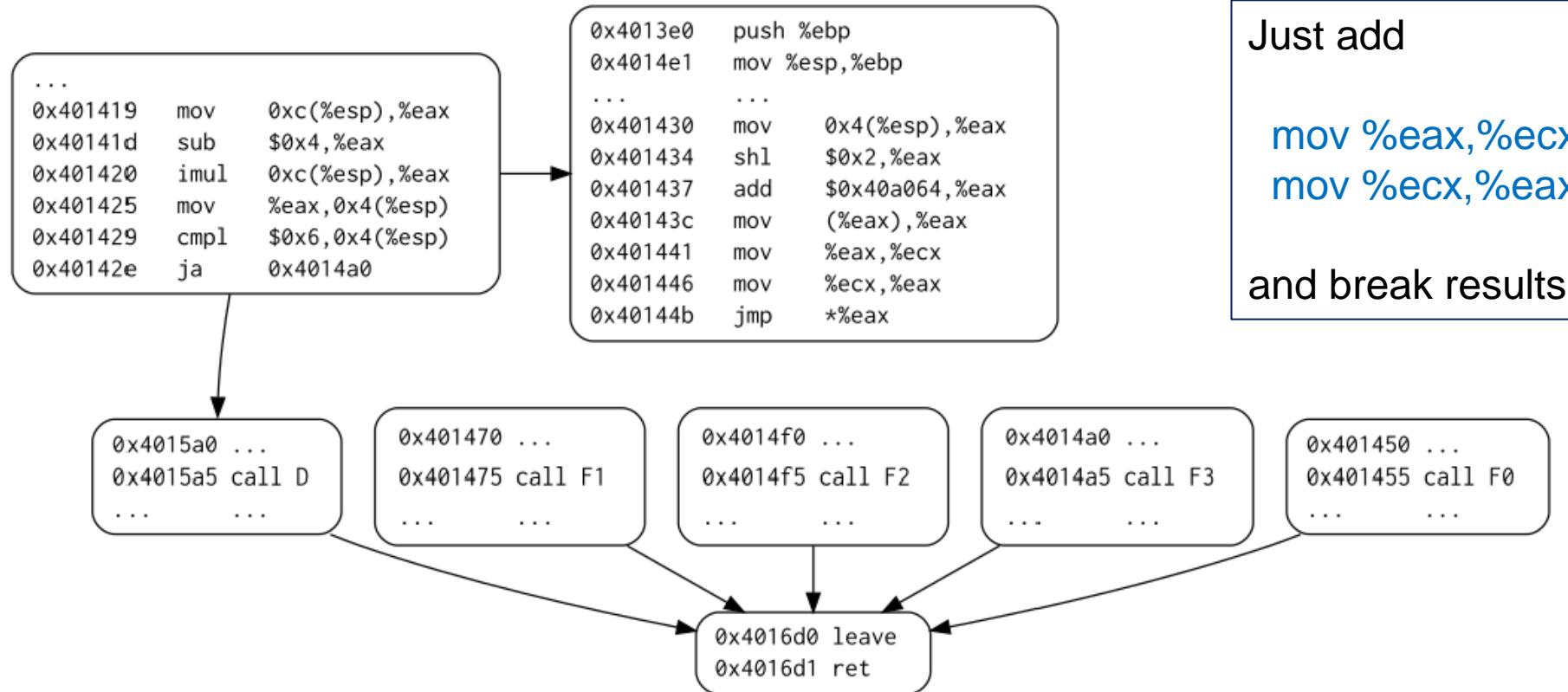
"val%d\n"

switch jump table

Assembly



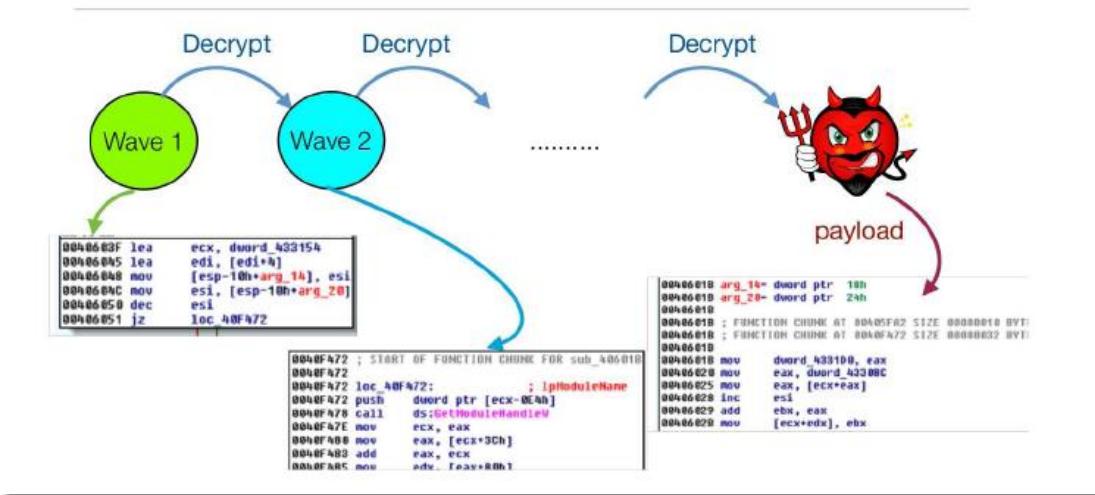
STATE-OF-THE-ART TOOLS ARE NOT ENOUGH



With IDA

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

[See later] CAN BECOME A NIGHTMARE WHEN OBFUSCATED



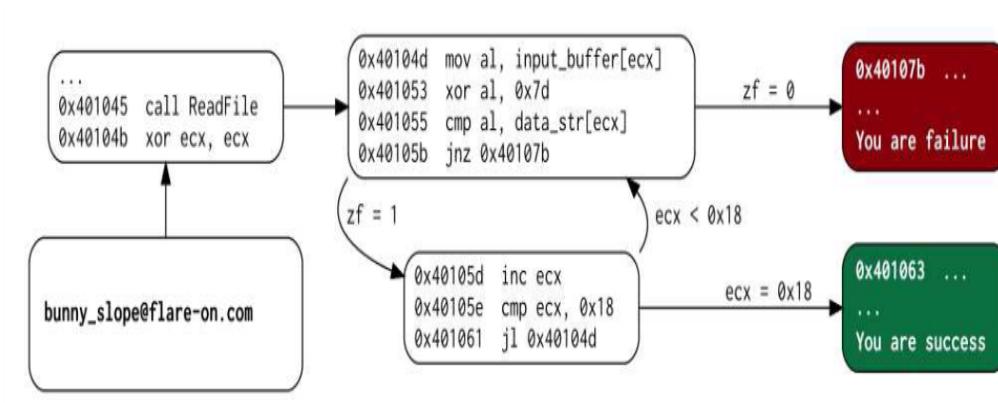
eg: $7y^2 - 1 \neq x^2$

(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>
```

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



EXAMPLE: VULNERABILITY DETECTION

Find vulnerabilities before the bad guys

- On the whole program
- At binary-level
- Know only the entry point and program input format

```
4800 6000 5dc3 5589 e5c7 0812 6000 00b8 4800 6000 5dc3 558
0000 00b8 4500 0000 5 0000 0000 0000 00b8 4500 000
bf0e 0821 0000 00b8 0000 0000 0000 0000 0000 0000 0000 000
e5c7 0540 bf0e 0822 0000 0000 0000 0000 0000 0000 0000 0000 000
5dc3 5589 e583 ec10 c705 00b8 4900 0000 0000 5dc3 5589 e583 ec1
0000 a148 bf0e 0883 f809 48bf 6e08 0100 0000 a148 bf0e 088
8b04 8548 e10b 08FF e0c6 0597 0002 0000 0000 8b04 8548 e10b 08F
00c6 45f9 00c6 45fa 00c7 45f7 00c6 45f8 00c6 45f9 00c6 45f
0000 00e9 d981 0000 c645 0548 bf0e 0882 0000 0000 d981 000
c645 f900 c645 fa01 807d f701 c645 f800 c645 f900 c645 fa0
48bf 6e08 0300 0000 807d fb00 750a c705 48bf 6e08 0300 000
fc00 750a c705 48bf 6e08 fb00 7410 807d fc00 750a c705 48b
fc00 7415 807d fb00 740f 0900 0000 807d fc00 7415 807d fb0
0600 0000 c988 0100 00e9 c705 48bf 6e08 0600 0000 c988 010
f701 c645 f800 c645 f900 0000 0000 c645 f701 c645 f800 c64
fc00 740f c705 48bf 6e08 c645 fa02 807d fc00 740f c705 48b
0100 00e9 5901 0000 c645 0400 0000 e95e 0100 00e9 5901 000
c645 f900 c645 fa03 807d f701 c645 f800 c645 f900 c645 fa0
fe00 750a c705 48bf 6e08 f000 7410 807d fe00 750a c705 48b
fc00 750a c705 48bf 6e08 0500 0000 807d fc00 750a c705 48b
fe00 740f c705 48bf 6e08 0300 0000 807d fc00 740f c705 48b
0100 0000 c645 0600 0000 e90e 0100 00e9 0901 000
c645 0043 f001 807d f701 c645 f800 c645 f901 c645 fa0
48bf 0400 0000 c9e4 807d f701 c645 f800 c645 f900 c645 fa0
0000 c645 f701 c645 f800 0000 00e9 df00 0000 c645 f701 c64
fa04 807d fc00 7410 807d c645 1000 c645 fa04 807d fc00 741
48bf 6e08 0700 0000 807d ff00 750a c705 48bf 6e08 0700 000
ff00 740f c705 48bf 6e08 fc00 7415 807d ff00 740f c705 48b
0000 00e9 9900 0000 c645 0600 0000 e99e 0000 00e9 9900 000
c645 f900 c645 fa05 807d f701 c645 f800 c645 f900 c645 fa0
fc00 7500 c705 48bf 6e08 f000 7410 807d fe00 750a c705 48b
fc00 750a c705 48bf 6e08 0800 0000 807d fc00 750a c705 48b
fe00 7506 807d ff00 740c 0900 0000 807d fe00 7506 807d ff0
0600 6000 eb4b eb49 c645 c705 48bf 6e08 0600 0000 eb4b eb4
c645 f901 c645 fa02 807d f701 c645 f800 c645 f901 c645 fa0
5dc3 5589 e5c7 0540 bf0e 00b8 5400 0000 5dc3 5589 e5c7 054
1800 6000 5dc3 5589 e5c7 0812 6000 00b8 4800 6000 5dc3 558
3000 00b8 4500 0000 5dc3 0540 bf0e 0820 0000 0000 0000 0000
bf0e 0821 0000 00b8 5800 5589 e5c7 0540 bf0e 082 USE 00b1
e5c7 0540 bf0e 0822 0000 0000 5dc3 5589 e5c7 054 082;
5dc3 5589 e583 ec10 c705 00b8 4900 0000 5dc3 5589 e583 ec1
3000 a148 bf0e 0883 f809 48bf 6e08 0100 0000 a148 bf0e 088
3b04 8548 e10b 08FF e0c6 0F87 0002 0000 8b04 8548 e10b 08F
00c6 45f9 00c6 45fa 00c7 45f7 00c6 45f8 00c6 45f9 00c6 45f
3000 00e9 d981 0000 c645 0548 bf0e 0882 0000 0000 d981 000
c645 f900 c645 fa01 807d f701 c645 f800 c645 f900 c645 fa0
48bf 6e08 0300 0000 807d fb00 750a c705 48bf 6e08 0300 000
fc00 750a c705 48bf 6e08 fb00 7410 807d fc00 750a c705 48b
fc00 7415 807d fb00 740f 0900 0000 807d fc00 7415 807d fb0
3000 0000 c988 0100 00e9 c705 48bf 6e08 0600 0000 c988 010
```

EXAMPLE: VULNERABILITY DETECTION

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)
firefox (CVE-2014-1512), chrome (CVE-2014-1713)

```
1 char *login , *passwords;
2 login=(char *) malloc( ... );
3 [...]
4 free( login ); // login is now a dangling pointer
5 [...]
6 passwords=(char *) malloc( ... ); // may re-allocate memory of *login
7 [...]
8 printf("%s\n" , login ); // security threat : may print the passwords !
```

CHALLENGE: In-depth exploration (example: use after free)

Find a needle in the heap !

- sequence of events, importance of aliasing
- strongly depend on implem of malloc and free

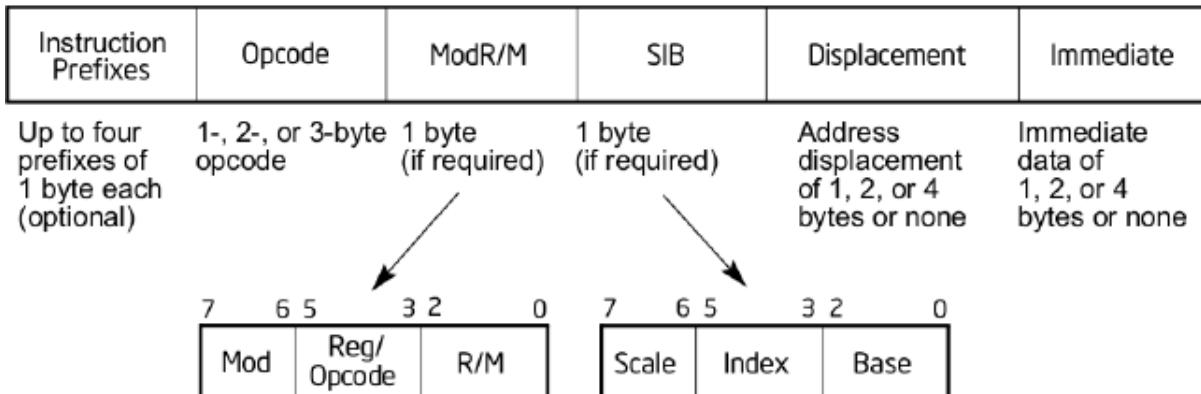
Dynamic: not enough

- Too incomplete



```
4800 6000 5dc3 5589 e5c7 0812 6000 00b8 4800 6000 5dc3 558  
0000 00b8 4500 0000 5dc3 5589 e5c7 0812 6000 00b8 4800 6000  
bf0e 0821 0000 00b8 5dc3 5589 e5c7 0812 6000 00b8 4500 000  
e5c7 0540 bf0e 0822 0000 00b8 5dc3 5589 e5c7 0540 bf0e 082  
5dc3 5589 e583 ec10 c705 00b8 4900 0000 5dc3 5589 e583 ec1  
0000 a148 bf0e 0883 f809 48bf 0e08 0100 0000 a148 bf0e 088  
8b04 8548 e10b 08FF e0c6 0f97 0002 0000 8b04 8548 e10b 08F  
00c6 45f9 00c6 45fa 00c7 45f7 00c6 45f8 00c6 45f9 00c6 45f  
0000 00e9 d981 0000 c645 0548 bf0e 0882 0000 00e9 d981 000  
c645 f900 c645 fa01 807d f701 c645 f800 c645 f900 c645 fa0  
48bf 0e08 0300 0000 807d f800 750a c705 48bf 0e08 0300 000  
fc00 750a c705 48bf 0e08 f800 7410 807d fc00 750a c705 48b  
fc00 7415 807d f800 740f 0900 0000 807d fc00 7415 807d fb0  
0600 0000 c988 0100 00e9 c705 48bf 0e08 0600 0000 c988 010  
f701 c645 f800 c645 f900 0000 0000 c645 f701 c645 f800 c64  
fc00 740f c705 48bf 0e08 c645 fa02 807d fc00 740f c705 48b  
0100 00e9 5901 0000 c645 0400 0000 e95e 0100 00e9 5901 000  
c645 f900 c645 fa03 807d f701 c645 f800 c645 f900 c645 fa0  
fe00 750a c705 48bf 0e08 f800 7410 807d fe00 750a c705 48b  
fc00 750a c705 48bf 0e08 0500 0000 807d fc00 750a c705 48b  
fe00 740f c705 48bf 0e08 0300 0000 807d fc00 740f c705 48b  
0100 0000 c645 0600 0000 e90e 0100 00e9 0901 000  
c645 0043 f804 807d f701 c645 f800 c645 f901 c645 fa0  
48bf 0400 0000 c9e4 f800 750f c705 48bf 0e08 0400 000  
0000 c645 f701 c645 f800 0000 00e9 df00 0000 c645 f701 c64  
fa04 807d fc00 7410 807d c645 1900 c645 fa04 807d fc00 741  
48bf 0e08 0700 0000 807d ff00 750a c705 48bf 0e08 0700 000  
ff00 740f c705 48bf 0e08 fc00 7415 807d ff00 740f c705 48b  
0000 00e9 9900 0000 c645 0600 0000 e99e 0000 00e9 9900 000  
c645 f900 c645 fa05 807d f701 c645 f800 c645 f900 c645 fa0  
fc00 750a c705 48bf 0e08 f800 7410 807d fe00 750a c705 48b  
fc00 750a c705 48bf 0e08 0800 0000 807d fc00 750a c705 48b  
fe00 7506 807d ff00 740c 0900 0000 807d fe00 7506 807d ff0  
0600 0000 eb4b eb49 c645 c705 48bf 0e08 0600 0000 eb4b eb4  
c645 f901 c645 fa02 807d f701 c645 f800 c645 f901 c645 fa0  
5dc3 5589 e5c7 0540 bf0e 00b8 5400 0000 5dc3 5589 e5c7 0541  
1800 0000 5dc3 5589 e5c7 0812 6000 00b8 4800 6000 5dc3 558!  
3000 00b8 4500 0000 5dc3 0540 bf0e 0820 0000 00b1 0001  
bf0e 0821 0000 00b8 5800 5589 e5c7 0540 bf0e 082 USE 00b1  
e5c7 0540 bf0e 0822 0000 0000 5dc3 5589 e5c7 0541 0821  
5dc3 5589 e583 ec10 c705 00b8 4900 0000 5dc3 5589 e583 ec1  
3000 a148 bf0e 0883 f809 48bf 0e08 0100 0000 a148 bf0e 088  
3b04 8548 e10b 08FF e0c6 0f87 0002 0000 8b04 8548 e10b 08F  
00c6 45f9 00c6 45fa 00c7 45f7 00c6 45f8 00c6 45f9 00c6 45f  
0000 00e9 d981 0000 c645 0548 bf0e 0882 0000 00e9 d981 000  
c645 f900 c645 fa01 807d f701 c645 f800 c645 f900 c645 fa0  
48bf 0e08 0300 0000 807d f800 750a c705 48bf 0e08 0300 000  
fc00 750a c705 48bf 0e08 f800 7410 807d fc00 750a c705 48b  
fc00 7415 807d f800 740f 0900 0000 807d fc00 7415 807d fb0  
0600 0000 c988 0100 00e9 c705 48bf 0e08 0600 0000 c988 0101
```

BONUS: (MULTI-)ARCHITECTURE SUPPORT



Example of x86

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

r0(/r)	AL	CL	DL	BL	AH	CH	DH
r16(/r)	AX	CX	DX	BX	SP	BP	SI
r32(/r)	EAX	ECX	EDX	EBX	ESP	EBP	ESI
mm(/r)	MM0	MM1	MM2	MM3	MM4	MM5	MM7
xmm(/r)	XMM0	XMM1	XMM2	XMM3	XMM4	XMM5	XMM7
sreg	ES	CS	SS	DS	FS	GS	res.
eee	CR0	invd	CR2	CR3	CR4	invd	invd
ooo	DR0	DR1	DR2	DR3	DR4 ¹	DR5 ¹	DR6
(In decimal) /digit (Opcode)	0	1	2	3	4	5	6
(In binary) REG =	000	001	010	011	100	101	110
Effective Address	Mod	R/M	Value of ModR/M Byte (in Hex)				
[EAX]	00	00	08	10	18	20	28
[ECX]	001	01	09	11	19	21	29
[EDX]	010	02	0A	12	1A	22	2A
[EBX]	011	03	0B	13	1B	23	2B
[sib]	100	04	0C	14	1C	24	2C
disp32	101	05	0D	15	1D	25	2D
[ESI]	110	06	0E	16	1E	26	2E
[EDI]	111	07	0F	17	1F	27	2F
[EAX]+disp8	81	008	40	48	50	58	60
[ECX]+disp8	001	41	49	51	59	61	69
[EDX]+disp8	010	42	4A	52	5A	62	7A
[EBX]+disp8	011	43	4B	53	5B	63	73
[sib]+disp8	100	44	4C	54	5C	64	74
[EBP]+disp8	101	45	4D	55	5D	65	75
[ESI]+disp8	110	46	4E	56	5C	66	76
[EDI]+disp8	111	47	4F	57	5F	67	77
[EAX]+disp32	10	000	80	88	90	98	A0
[ECX]+disp32	001	81	89	91	99	A1	A9
[EDX]+disp32	010	82	8A	92	9A	A2	A8
[EBX]+disp32	011	83	8B	93	9B	A3	A8
[sib]+disp32	100	84	8C	94	9C	A4	A4
[EBP]+disp32	101	85	8D	95	9D	A5	A5
[ESI]+disp32	110	86	8E	96	9E	A6	A6
[EDI]+disp32	111	87	8F	97	9F	A7	B7
AL/AX/EAX/ST0/MM0/XMM0	11	000	C0	C8	D0	D8	E0
CL/CX/ECX/ST1/MM1/XMM1	001	C1	C9	D1	D9	E1	E9
DL/DX/EDX/ST2/MM2/XMM2	010	C2	CA	D2	DA	E2	FA
BL/BX/EBX/ST3/MM3/XMM3	011	C3	CB	D3	DB	E3	F3
AH/SP/ESP/ST4/MM4/XMM4	100	C4	CC	D4	DC	E4	F4
CH/RP/FRP/ST5/MM5/XMM5	101	C5	CD	D5	DD	E5	FD

THE SITUATION

- **Binary-level security analysis is necessary**
- **Binary-level security analysis is highly challenging (*)**
- **Standard tools are not enough – experts need better help!**

(*) i.e., more challenging than source code analysis

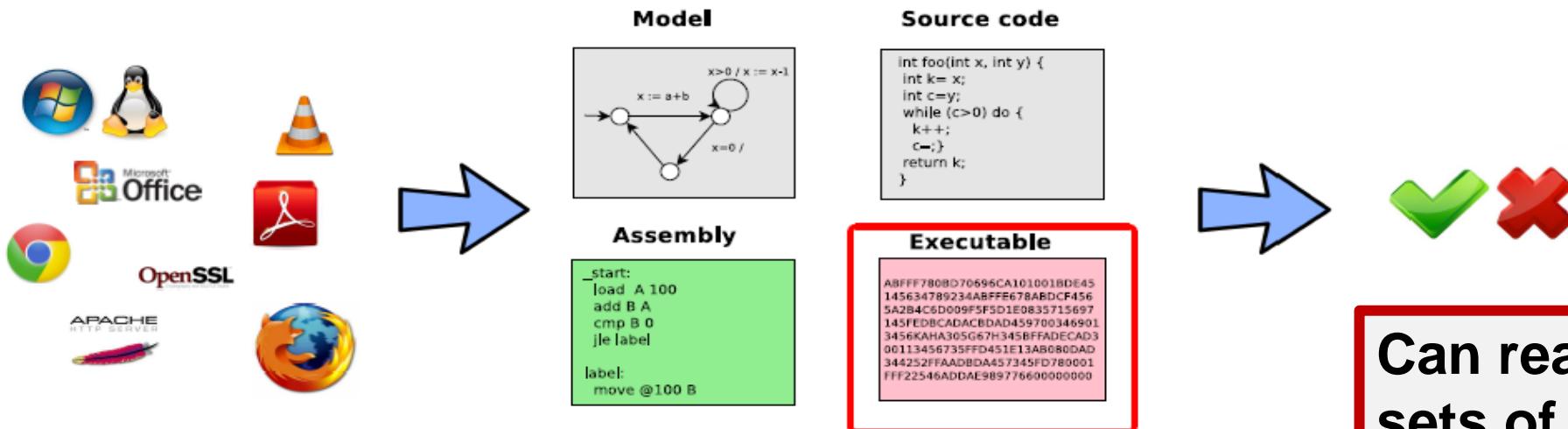
- **Static (syntactic): too fragile**
- **Dynamic: too incomplete**

SOLUTION? BINARY-LEVEL SEMANTIC ANALYSIS

Semantic tools help make sense of binary

- Develop the next generation of binary-level tools !
- motto : leverage formal methods from safety critical systems

**Semantic preserved
by compilation or
obfuscation**



**Can reason about
sets of executions**

Advantages

- more robust than syntactic
- more thorough than dynamic

Challenges

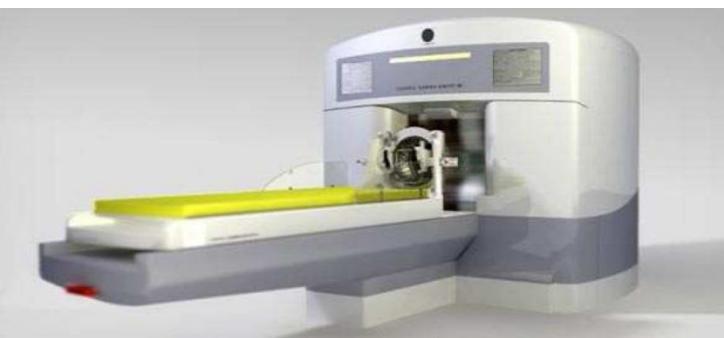
- source-level \mapsto binary-level
- safety \mapsto security
- many (complex) architectures

- Why binary-level analysis?
- Some background on source-level formal methods
- The hard journey from source to binary
- A few case-studies
- Conclusion

BACK IN TIME: THE SOFTWARE CRISIS (1969)

The major cause of the software crisis is that the machines have become several orders of magnitude more powerful! To put it quite bluntly : as long as there were no machines, programming was no problem at all; when we had a few weak computers, programming became a mild problem, and now we have gigantic computers, programming has become an equally gigantic problem.

- Edsger Dijkstra, The Humble Programmer (EWD340)



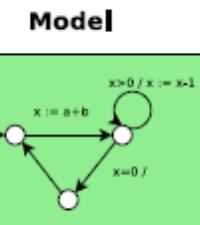
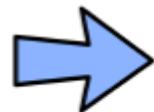
http://en.wikipedia.org/wiki/List_of_software_bugs

Testing can only reveal the presence of errors but never their absence.

- E. W. Dijkstra (Notes on Structured Programming, 1972)

ABOUT FORMAL METHODS

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



Source code

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



Success in safety-critical



Key concepts : $M \models \varphi$

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

Kind of properties

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

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

A DREAM COME TRUE ... IN CERTAIN DOMAINS (2)

Ex : Microsoft

Verification of drivers [SDV]

- conformance to MS driver policy
- home developers
- and third-party developers



Things like even software verification, this has been the Holy Grail of computer science for many decades but now in some very key areas, for example, driver verification we're building tools that can do actual proof about the software and how it works in order to guarantee the reliability.

- Bill Gates (2002)

Semantics

- Precise meaning for the domain of evaluation and the effect of instructions
- Operational semantics = « interpreter »

Properties

- From Invariants / reachability to safety/liveness/hyper-properties/...
- On software: mostly invariants and reachability

Algorithms:

- Historically: Weakest precondition, Abstract interpretation, model checking
- Correctness: the analysis explores only behaviors of interest
- Completeness: the analysis explores at least all behaviors of interest

Trends:

- Frontier between techniques disappear
- master abstraction (correct xor complete)
- reduction to logic
- sweet spots

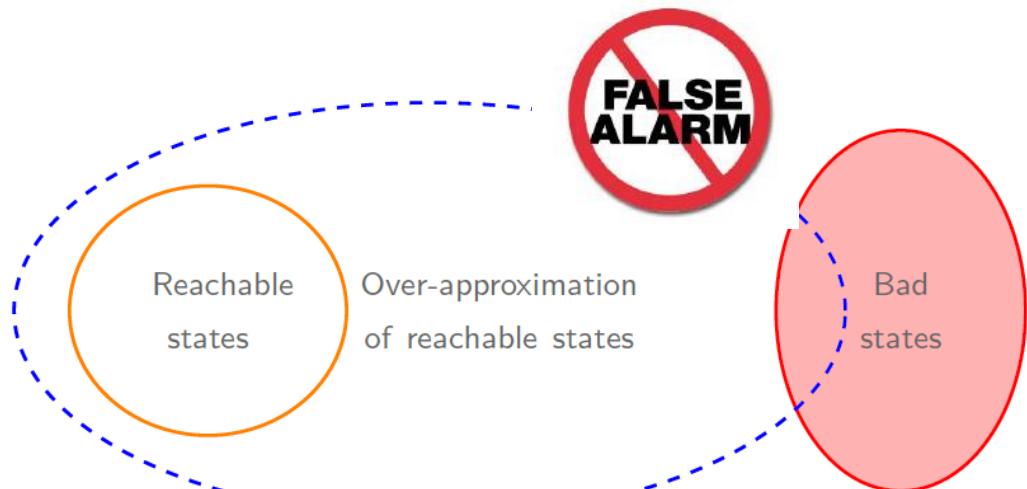
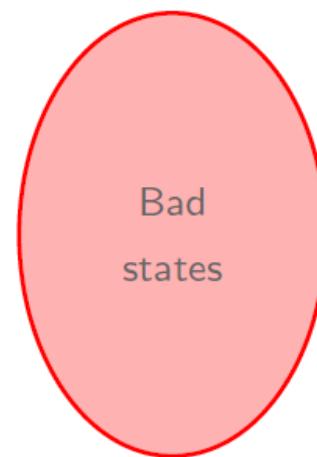
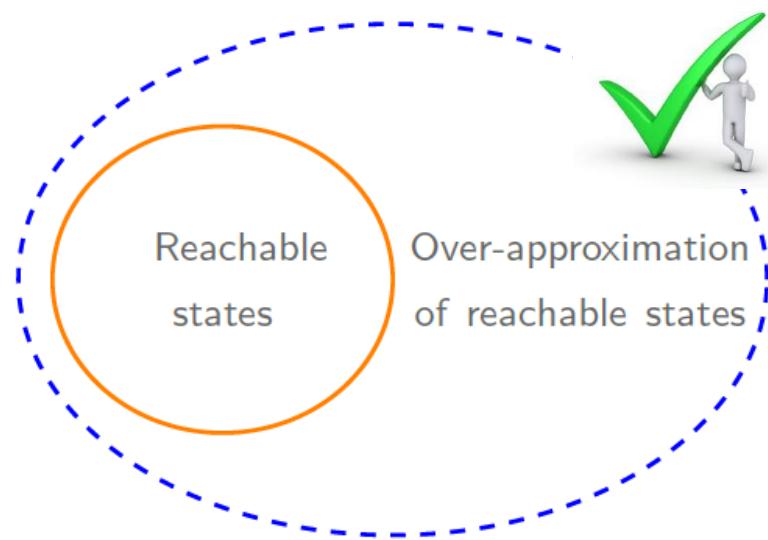
Next:

- AI: complete (can prove invariants) -- 1977
- DSE: correct (can find bugs) -- 2005

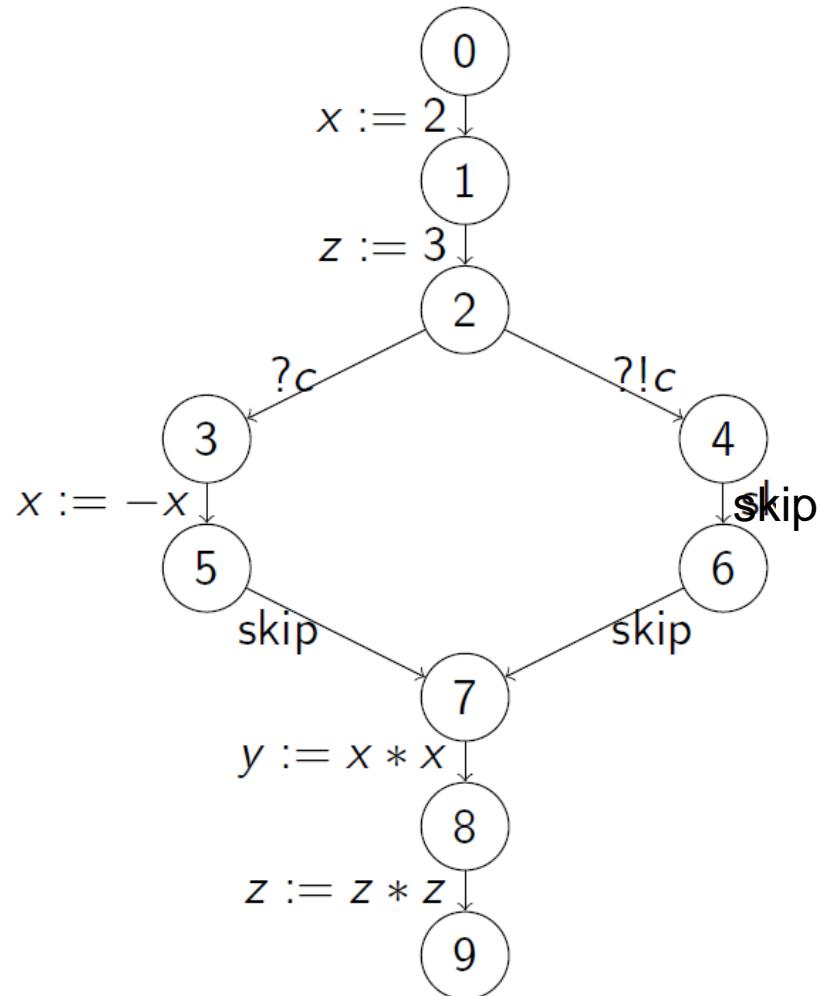
- Representative
- Industrial successes at source-level
- Adaptation to binary: very different situations

ABSTRACT INTERPRETATION

$$(\mathcal{P}(\text{states}), \cup, \cap, \rightarrow) \xrightleftharpoons[\alpha]{\gamma} (\text{states}^\#, \sqcup, \sqcap, \rightarrow^\#)$$



ABSTRACT INTERPRETATION IN PRACTICE



noeud	c	x	y	z
0	T	T	T	T
1	T	2	T	T
2	T	2	T	3
ip	T	2	T	3
4	0	2	T	3
5	T	-2	T	3
6	0	2	T	3
7	T	T2	T	3
8	T	T	T	3
9	T	T	T	9

Key points:

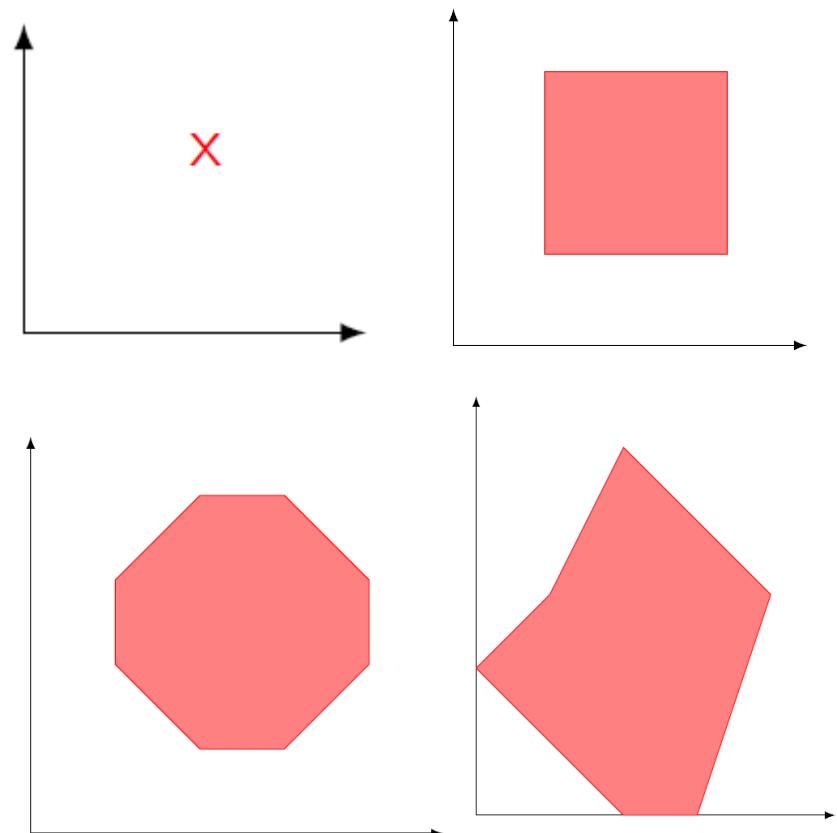
- Infinite data: abstract domain
- Path explosion: merge
- Loops: widening

In practice:

- Tradeoff between cost and precision
- Tradeoff between generic & dedicated domains

It is sometimes simple and useful

- taint, pointer nullness, typing

Big successes: Astrée, Frama-C, Clousot

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;  
}
```

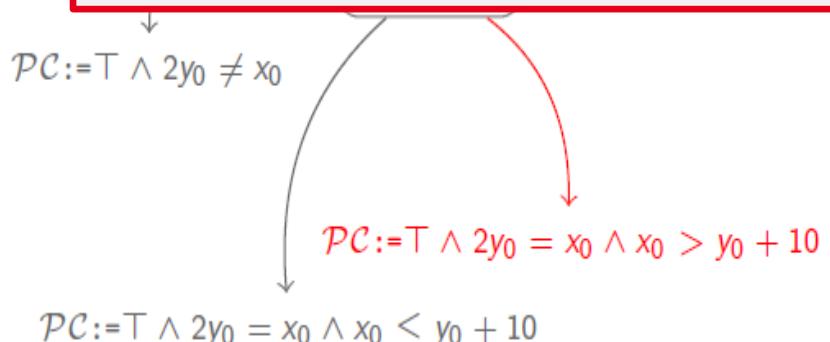
- given a path of the program
- automatically find input that follows the path
- then, iterate over all paths

$$\begin{array}{l}\sigma := \emptyset \\ \mathcal{PC} := T\end{array}$$

Perfect for intensive testing

- **Correct, relatively complete**
- **No false alarm**
- **Robust**
- **Scale in some ways**

// incomplete



DSE: PATH PREDICATE COMPUTATION (DSE, Godefroid 2005)

Loc	Instruction
0	input(y,z)
1	w := y+1
2	x := w + 3
3	if (x < 2 * z) (branche True)
4	if (x < z) (branche False)

let $W_1 \triangleq Y_0 + 1$ in
let $X_2 \triangleq W_1 + 3$ in
 $X_2 < 2 \times Z_0 \wedge X_2 \geq Z_0$

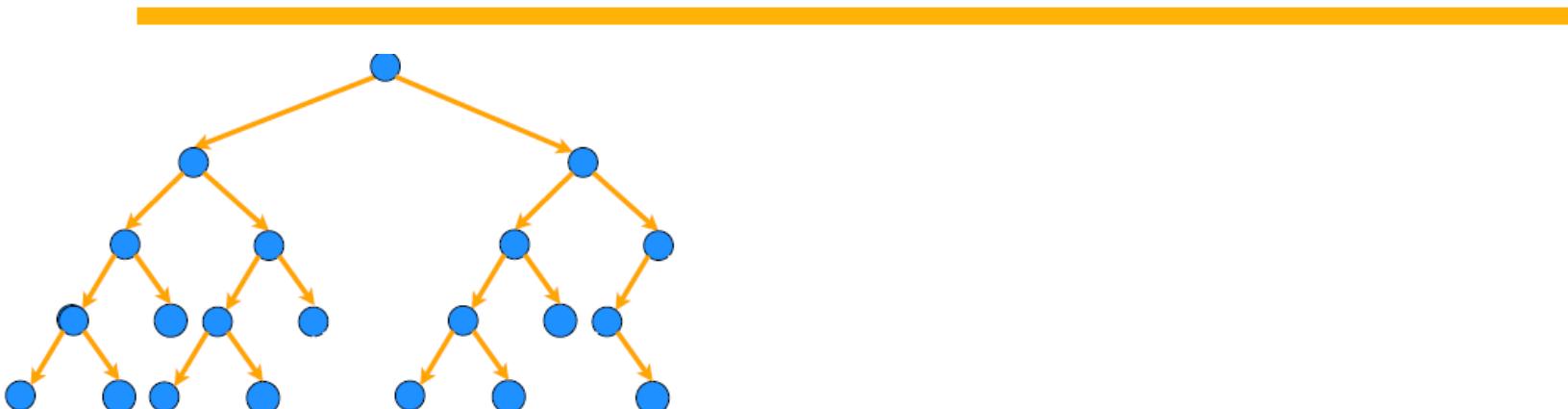
DSE: GLOBAL PROCEDURE

(DSE, Godefroid 2005)

input : a program P

output : a test suite TS covering all feasible paths of $Paths^{\leq k}(P)$

- pick a path $\sigma \in Paths^{\leq k}(P)$
- compute a *path predicate* φ_σ of σ
- solve φ_σ for satisfiability
- SAT(s) ? get a new pair $< s, \sigma >$
- loop until no more path to cover



ABOUT ROBUSTNESS (imo, the major advantage)

Goal = find input leading to ERROR
(assume we have only a solver for linear integer arith.)

```
g(int x) {return x*x; }  
f(int x, int y) {z=g(x); if (y == z) ERROR; else OK }
```

Symbolic Execution

- create a subformula $z = x * x$, out of theory [FAIL]

Dynamic Symbolic Execution

- first concrete execution with $x=3$, $y=5$ [goto OK]
- during path predicate computation, $x * x$ not supported
 - . x is concretized to 3 and z is forced to 9
- resulting path predicate : $x = 3 \wedge z = 9 \wedge y = z$
- a solution is found : $x=3$, $y=9$ [goto ERROR] [SUCCESS]

« concretization »

- Keep going when symbolic reasoning fails
- Tune the tradeoff genericity - cost

Three key ingredients

- Path predicate & solving
- Path enumeration
- C/S policy

Limits

- #paths -> better heuristics (?), state merging, distributed search, path pruning, adaptation to coverage objectives, etc.
- solving cost -> preprocessing, caching, incremental solving, aggressive concretization (good?)
[wait for better solvers ☺]
- Preconditions/postconditions/advanced stubs

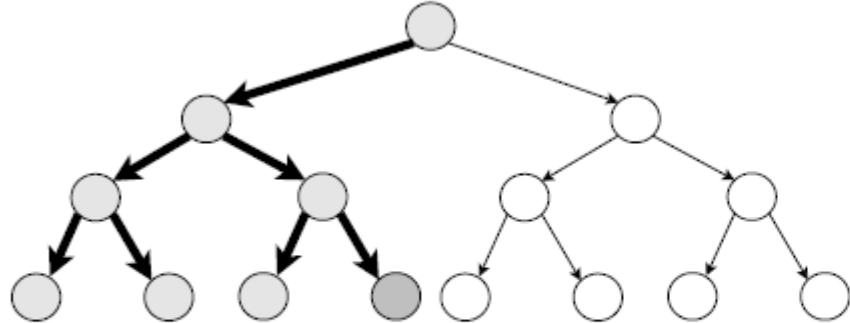
DSE: PATH PREDICATE MAY BE COMPLICATED

$x := a + b$

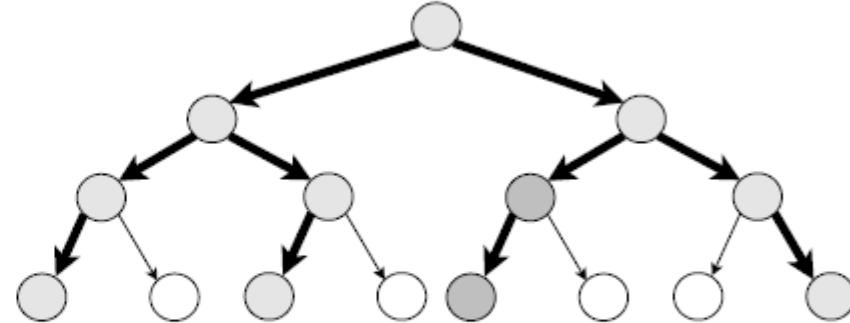
$X_{n+1} = A_n + B_n$

$store(M, \text{addr}(X), \text{load}(M, \text{addr}(A)) + \text{load}(M, \text{addr}(B)))$

```
let tmpA = load(M,addr(A)) @ load(M,addr(A)+1) @ load(M,addr(A)+2)
and tmpB = load(M,addr(B)) @ load(M,addr(B)+1) @ load(M,addr(B)+2)
in
let nX = tmpA+tmpB
in
  store(
    store(
      store(
        store(M,addr(X), nX[0]),
        addr(X) + 1, nX[1]),
        addr(X) + 2, nX[2]))
```

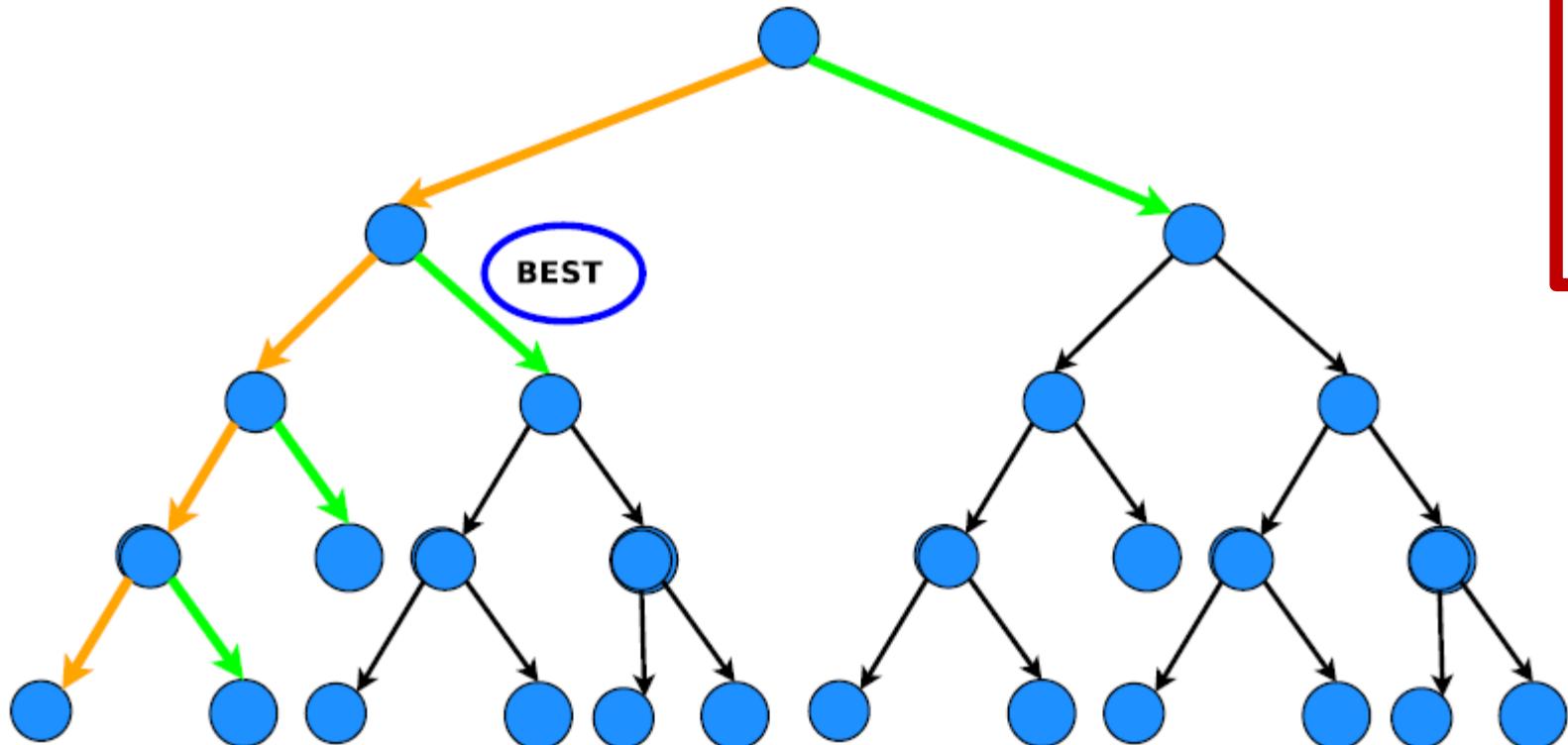


(a) DFS



(b) BFS

- **Search heuristics matters**
- But no good choice (hint: DFS is often the worst)
- The engine must provide flexibility



Generic engine

- Score each active prefix
 - Pick the best & expand
 - Easy encoding of many heuristics

Robustness : what if the instruction cannot be reasoned about ?

- missing code, self-modification
- hash functions, dynamic memory accesses, NLA operators

program	path predicate	concretization	symbolization
<pre>input: a, b x := a × b x := x + 1 //assert x > 10</pre>	$\begin{aligned} \varphi_1 = & x_1 = a \times b \\ \wedge & x_2 = x_1 + 1 \\ \wedge & x_2 > 10 \end{aligned}$ (φ_1)	$\begin{aligned} a &= 5 \\ \wedge & x_1 = 5 \times b \\ \wedge & x_2 = x_1 + 1 \\ \wedge & x_2 > 10 \end{aligned}$ (φ_2)	$\begin{aligned} x_1 &= \text{fresh} \\ \wedge & x_2 = x_1 + 1 \\ \wedge & x_2 > 10 \end{aligned}$ (φ_3)

Solutions

- **Concretization** : replace by runtime value [lose completeness]
- **Symbolization** : replace by fresh variable [lose correctness]

C/S POLICIES (2)

Consider the following situation

- instruction $x := @(\mathbf{a} * \mathbf{b})$
 - your tool documentation says : “*memory accesses are concretized*”
 - suppose that at runtime : $\mathbf{a} = 7, \mathbf{b} = 3$
-

What is the intended meaning ? [perfect reasoning : $x == \text{select}(M, a \times b)$]

CS1 : $x == \text{select}(M, 21)$ [incorrect]

CS2 : $x == \text{select}(M, 21) \wedge a \times b == 21$ [minimal]

CS3 : $x == \text{select}(M, 21) \wedge a == 7 \wedge b == 3$ [atomic]

No best choice, depends on the context

- acceptable loss of correctness / completeness ?
- is it mandatory to get rid off \times ?

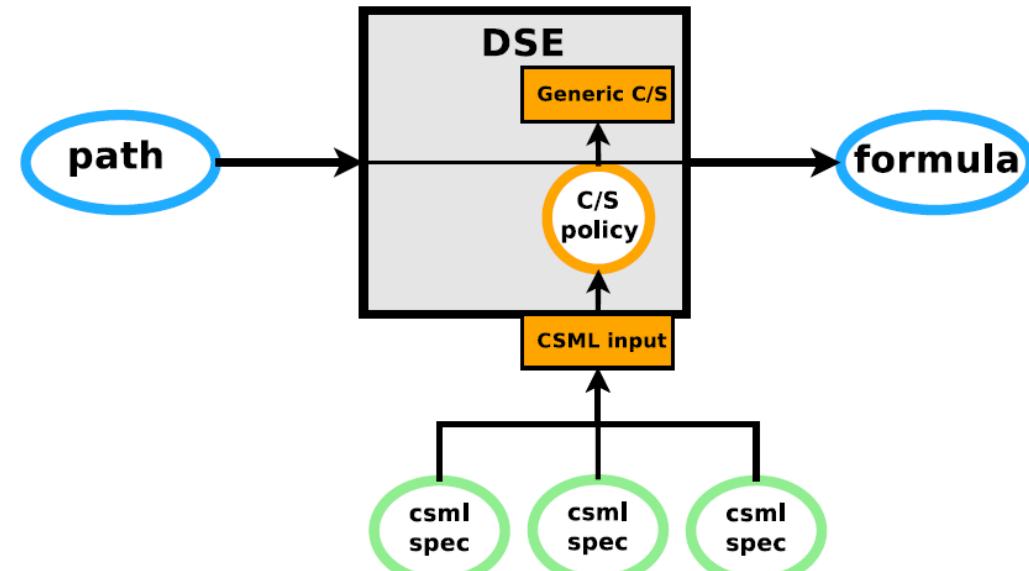
- **C/S policy matters**
- But no good choice
- The engine must provide flexibility

C/S POLICIES (3)



Generic engine

- C/S specification
- DSE parametrized by C/S

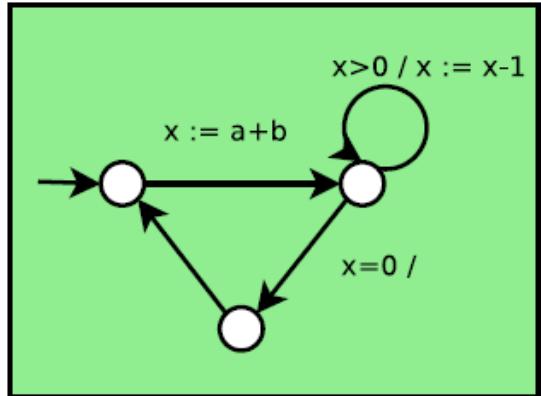


$$cs : \text{loc} \times \text{instr} \times \text{state} \times \text{expr} \mapsto \left\{ \begin{array}{ll} \mathcal{C} & \text{concretization} \\ \mathcal{S} & \text{symbolization} \\ \mathcal{P} & \text{propagation} \end{array} \right\}$$

- Why binary-level analysis?
- Some background on source-level formal methods
- The hard journey from source to binary
- A few case-studies
- Conclusion

NOW: BINARY-LEVEL SECURITY

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  
344252FFAABDBA457345FD780001  
FFF22546ADDAE9897766000000000
```

THE HARD JOURNEY FROM SOURCE TO BINARY

Low-level semantics of data

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

Wanted

- robustness
- precision
- scale

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

DSE is quite easy to adapt

- thx to SMT solvers (arrays+bitvectors)
- thx to concretization
- yet, performance degrades

Problems

- Low-level control: jump eax
- Low-level data: memory
- Low-level data: flags

AI is much more complicated

- Even for « normal » code
- btw, cannot expect better than source-level precision

Problem solved: multi-architecture

- rely on some IR

FULL DISCLOSURE: the BINSEC tool

Semantic analysis for binary-level security

- Help make sense of binary
- more robust than syntactic
- more exhaustive than dynamic

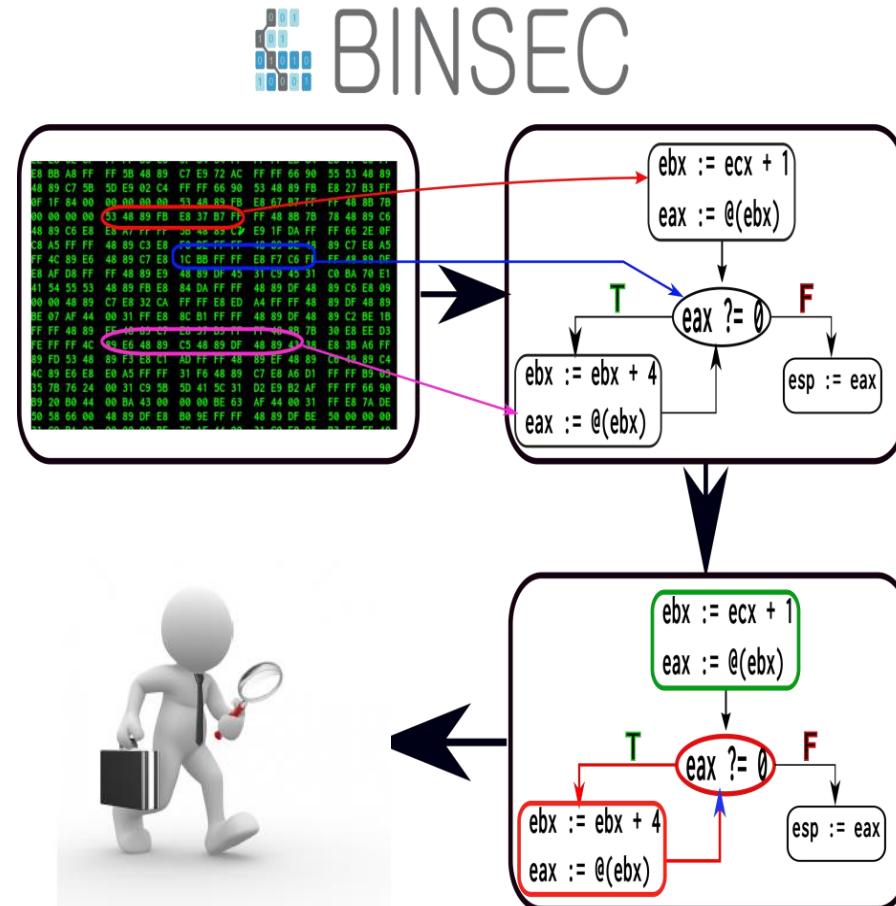
Some features

- Help to recover a simple model
- Identify feasible events (+ input)
- Identify infeasible events (eg, protections)
- Multi-architecture

Challenges

- Binary analysis
- **Scalability**
- **Robustness** w.r.t obfuscation

Still very young!



UNDER THE HOOD

x86

```
ABFFF780BD70696CA101001BDE45
145634789234ABFFE678ABDCF456
5A2B4C6D009F5F5D1E0835715697
145FEDBCADACBDAD459700346901
3456KAHA305G67H345BFFADECAD3
00113456735FFD451E13AB080DAD
344252FFAADD8DA457345FD780001
FFF22546ADDAE9897766000000000
```

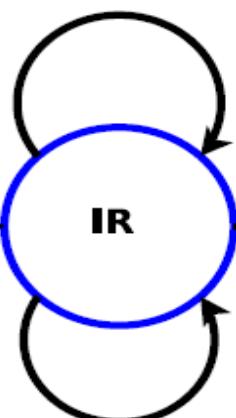
ARM

```
ABFFF780BD70696CA101001BDE45
145634789234ABFFE678ABDCF456
5A2B4C6D009F5F5D1E0835715697
145FEDBCADACBDAD459700346901
3456KAHA305G67H345BFFADECAD3
00113456735FFD451E13AB080DAD
344252FFAADD8DA457345FD780001
FFF22546ADDAE9897766000000000
```

...

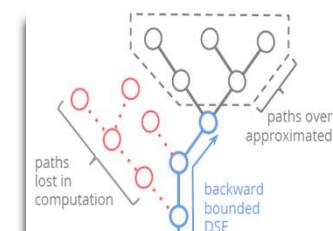
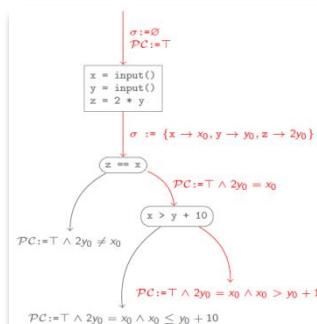
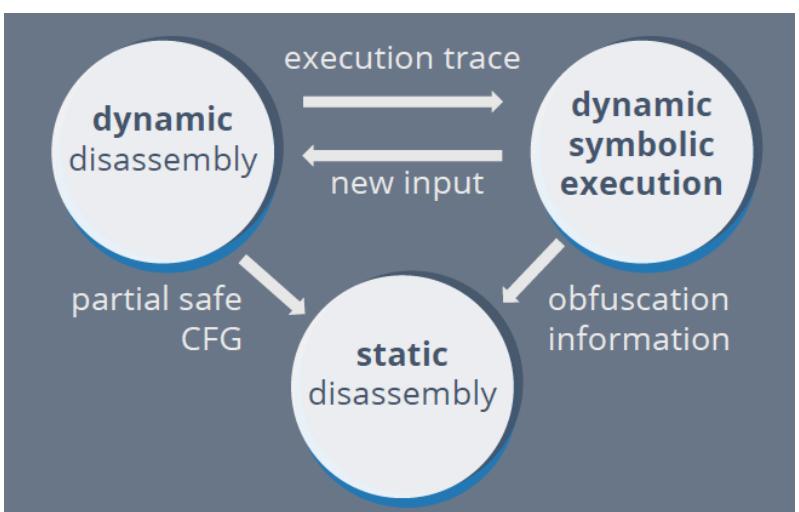
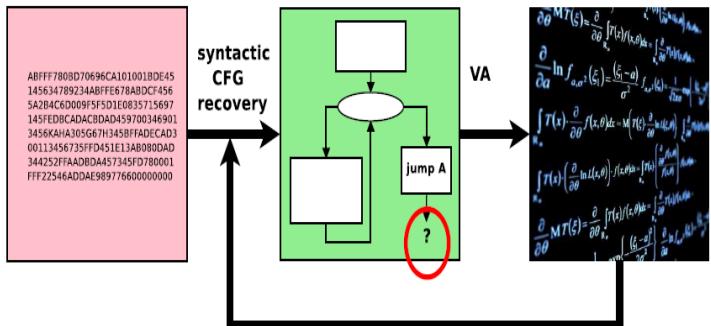
```
ABFFF780BD70696CA101001BDE45
145634789234ABFFE678ABDCF456
5A2B4C6D009F5F5D1E0835715697
145FEDBCADACBDAD459700346901
3456KAHA305G67H345BFFADECAD3
00113456735FFD451E13AB080DAD
344252FFAADD8DA457345FD780001
FFF22546ADDAE9897766000000000
```

Static analysis

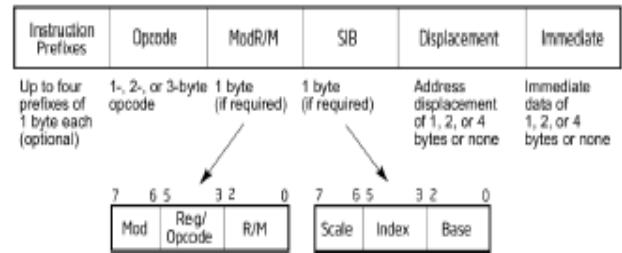


Symbolic execution

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



INTERMEDIATE REPRESENTATION



81 c3 57 1d 00 00 ^{x86reference} ⇒ ADD EBX 1d57

- `lhs := rhs`
- `goto addr, goto expr`
- `ite(cond)? goto addr`

- Concise
- Well-defined
- Clear, side-effect free

```
(0x29e,0) tmp := EBX + 7511;
(0x29e,1) OF := (EBX{31,31}=7511{31,31}) && (EBX{31,31}<>tmp{31,31});
(0x29e,2) SF := tmp{31,31};
(0x29e,3) ZF := (tmp = 0);
(0x28e,4) AF := ((extu (EBX{0,7}) 9) + (extu 7511{0,7} 9)){8,8};
(0x29e,6) CF := ((extu EBX 33) + (extu 7511 33)){32,32};
(0x29e,7) EBX := tmp; goto (0x2a4,0)
```

INTERMEDIATE REPRESENTATION + simplifications

program	native loc	DBA loc	opt (DBA)		
			time	loc	red
bash	166K	559K	673.61s	389K	30.45%
cat	8K	23K	18.54s	18K	23.02%
echo	4K	10K	6.96s	8K	24.26%
less	23K	80K	69.99s	55K	30.96%
ls	19K	63K	65.69s	44K	30.58%
mkdir	8K	24K	19.74s	17K	29.50%
netstat	17K	50K	52.59s	40K	20.05%
ps	12K	36K	36.99s	27K	23.98%
pwd	4K	11K	7.69s	9K	23.56%
rm	10K	30K	24.93s	22K	25.24%
sed	10K	32K	28.85s	23K	26.20%
tar	64K	213K	242.96s	154K	27.48%
touch	8K	26K	24.28s	18K	27.88%
uname	3K	10K	6.99s	8K	23.62%

Approach

- Inspired from standard compiler optim
- Targets : flags & temp
- Sound : w.r.t. incomplete CFG
- Inter-procedural (summaries)

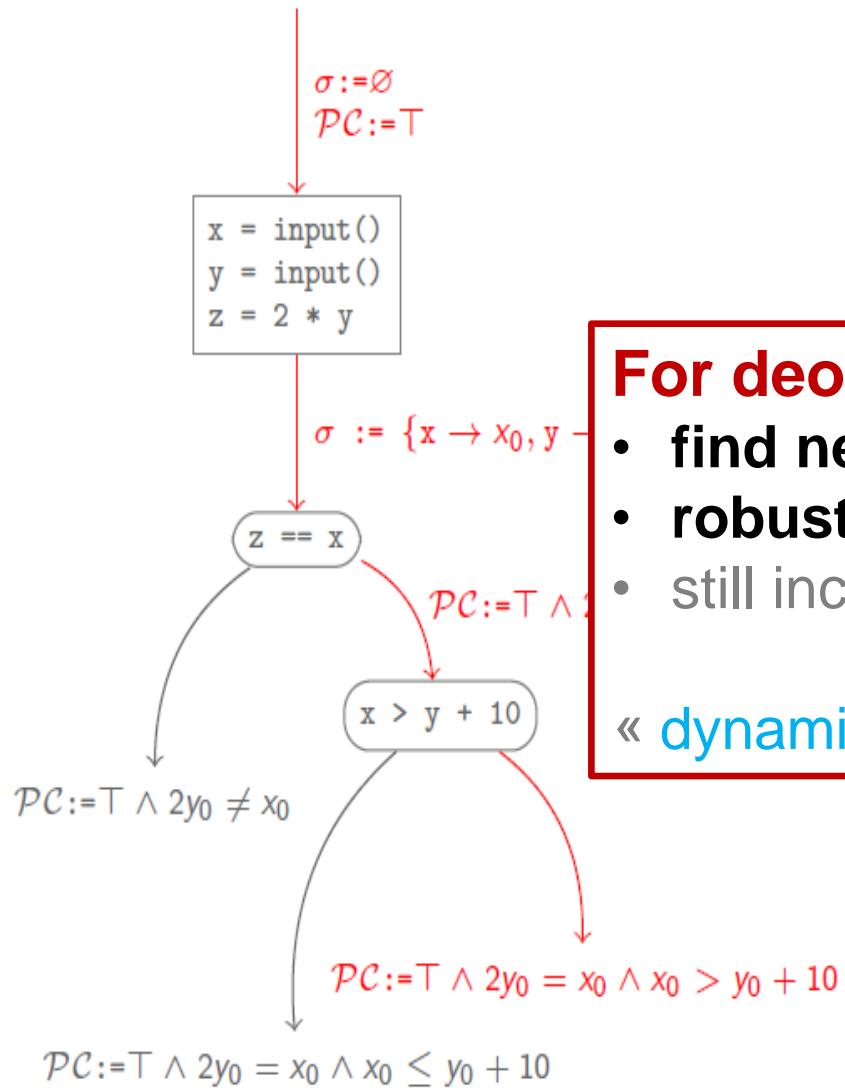
- **IR level**
- **machine-instruction level**
- **program level**

	reduction			
	time	dba instr	tmp assigns	flag assigns
BINSEC	1279.81s	28.64%	88.00%	75.04%

BINARY-LEVEL DSE (Godefroid)

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

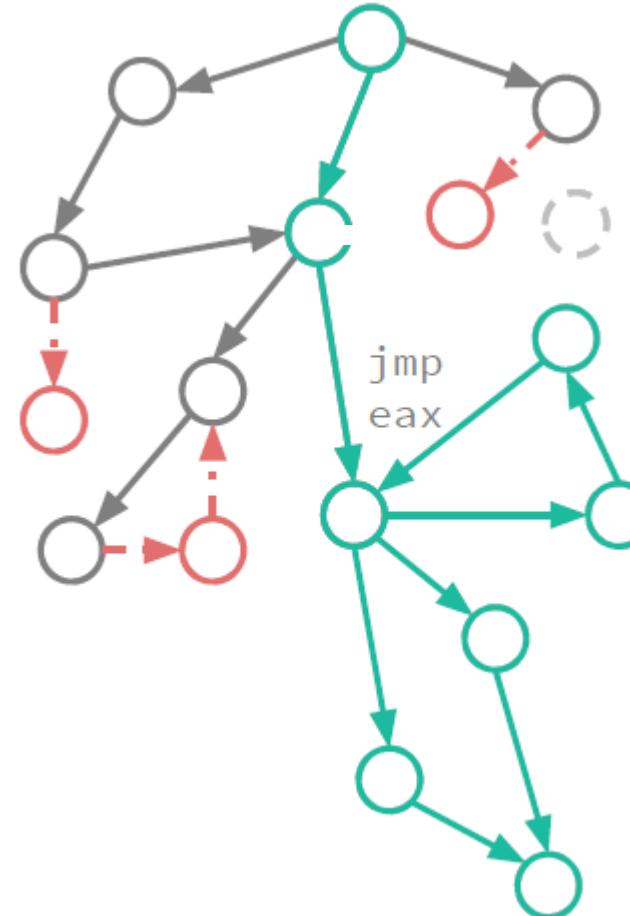
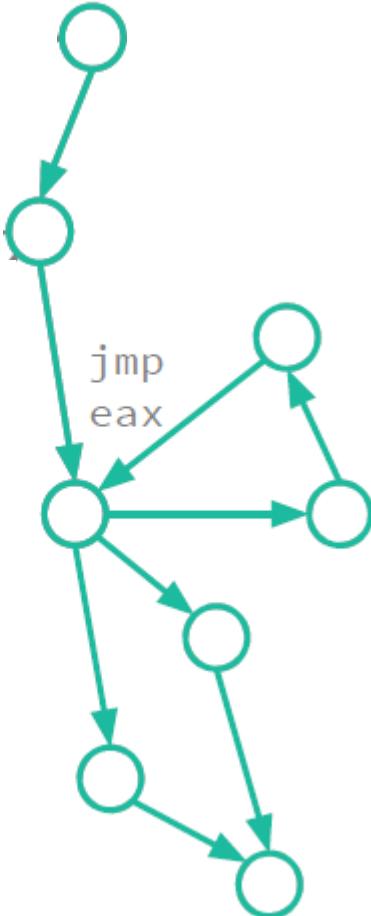
- given a path of the program
- automatically find input that follows the path
- then, iterate over all paths

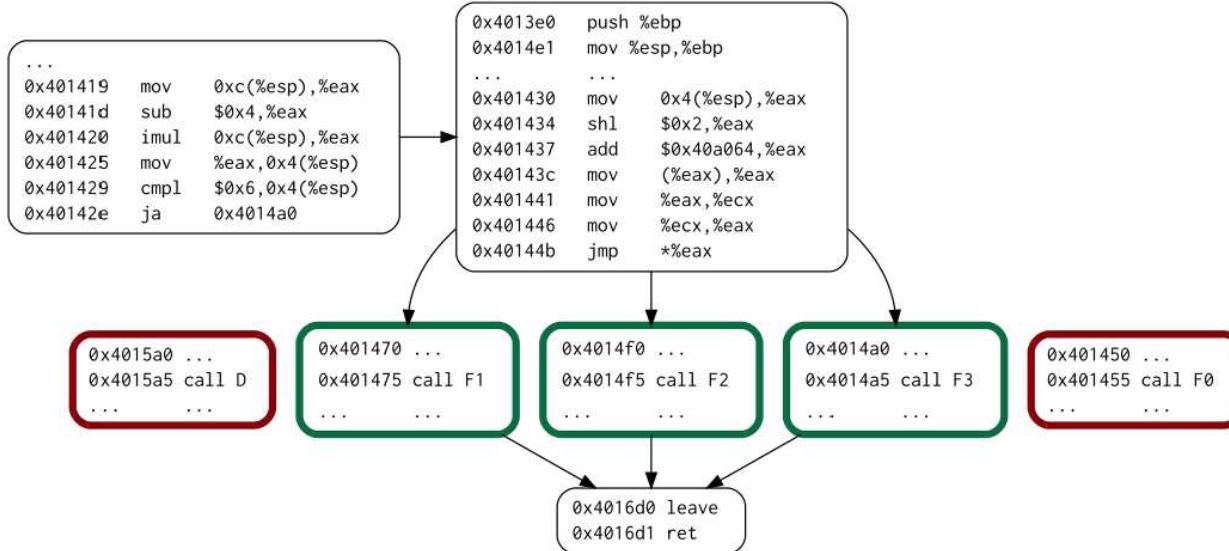


- For deobfuscation**
- find new real paths
 - robust
 - still incomplete

« dynamic analysis on steroids »

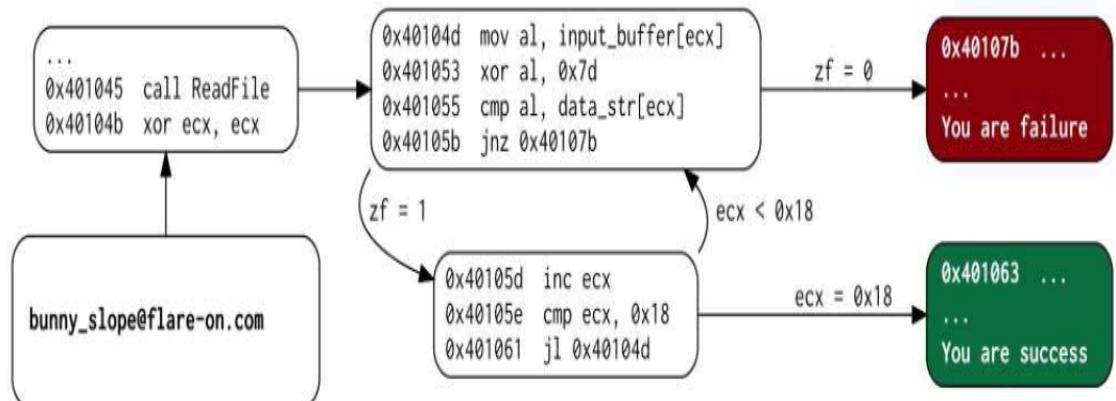
DSE COMPLEMENTS DYNAMIC ANALYSIS



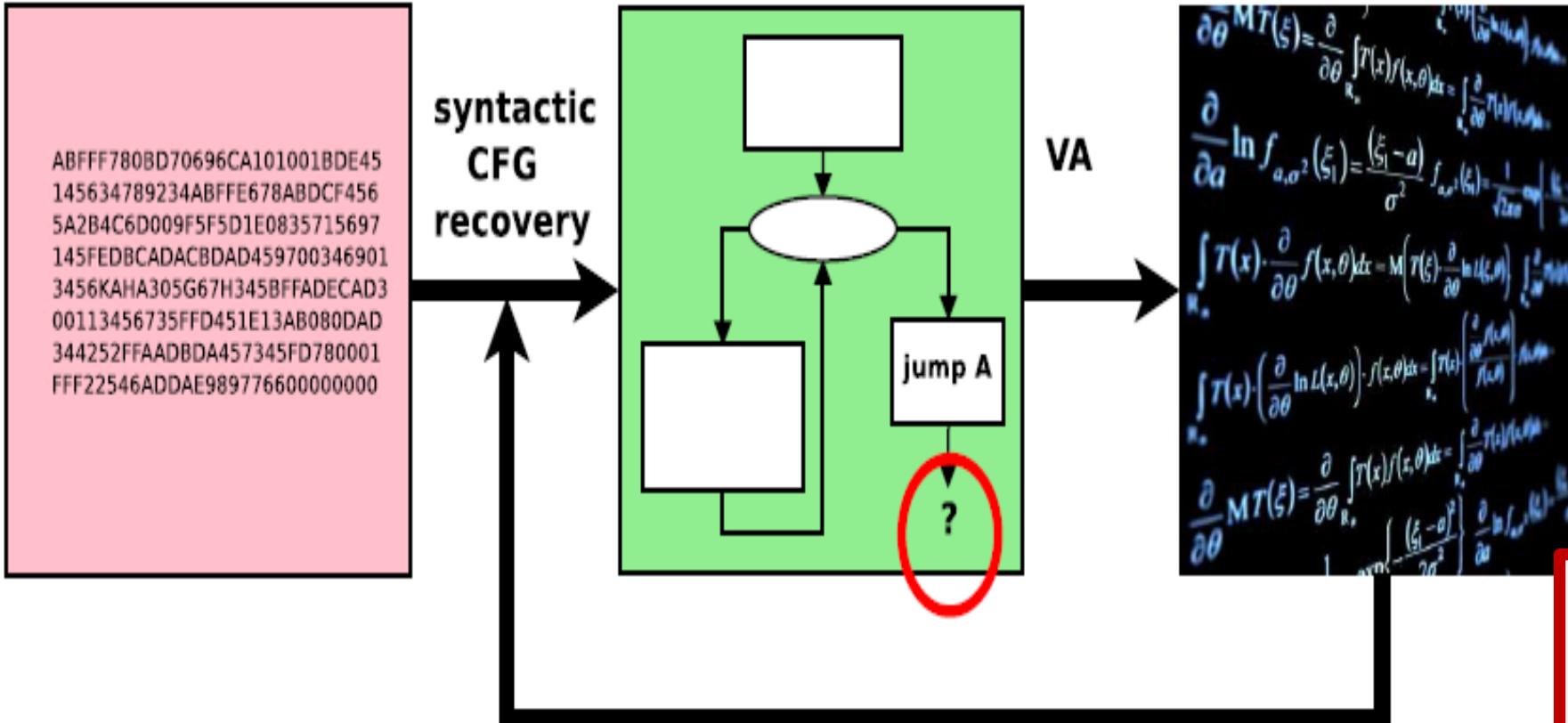


With IDA + BINSEC

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



ABSTRACT INTERPRETATION IS VERY VERY HARD ON BINARY CODE



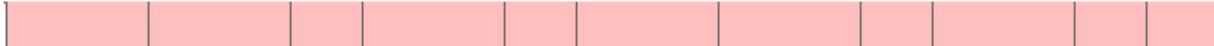
- Problems**
- Jump eax
 - memory
 - Bit reasoning

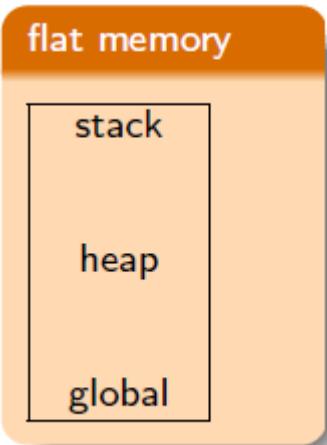
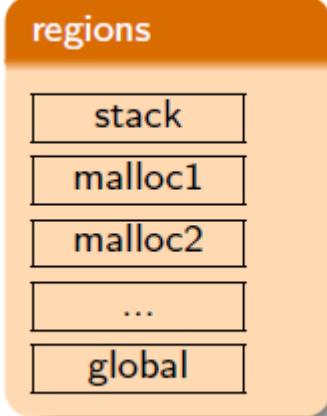
ISSUE: GLOBAL MEMORY

$\text{@}(r, 0) := \dots$ 

$\text{@}(r, T) := \dots$ 

$\text{@}(7) := \dots$ 

$\text{@}(T) := \dots$ 



- Problems
- Jump eax
 - memory
 - Bit reasoning

ISSUE: LACK of HIGH-LEVEL STRUCTURE

- Problems**
- Jump eax
 - memory
 - Bit reasoning

High-level conditions translated into low-level flag predicates

```
if (ax > bx) X = -1;  
else X = 1;
```

```
OF := ((ax{31,31}≠bx{31,31}) &  
       (ax{31,31}≠(ax-bx){31,31}));  
SF := (ax-bx) < 0;  
ZF := (ax-bx) = 0;  
if (¬ ZF ∧ (OF = SF)) goto 11  
X := 1  
goto 12  
11: X := -1  
12:
```

Condition on flags, not on register (nor stack)

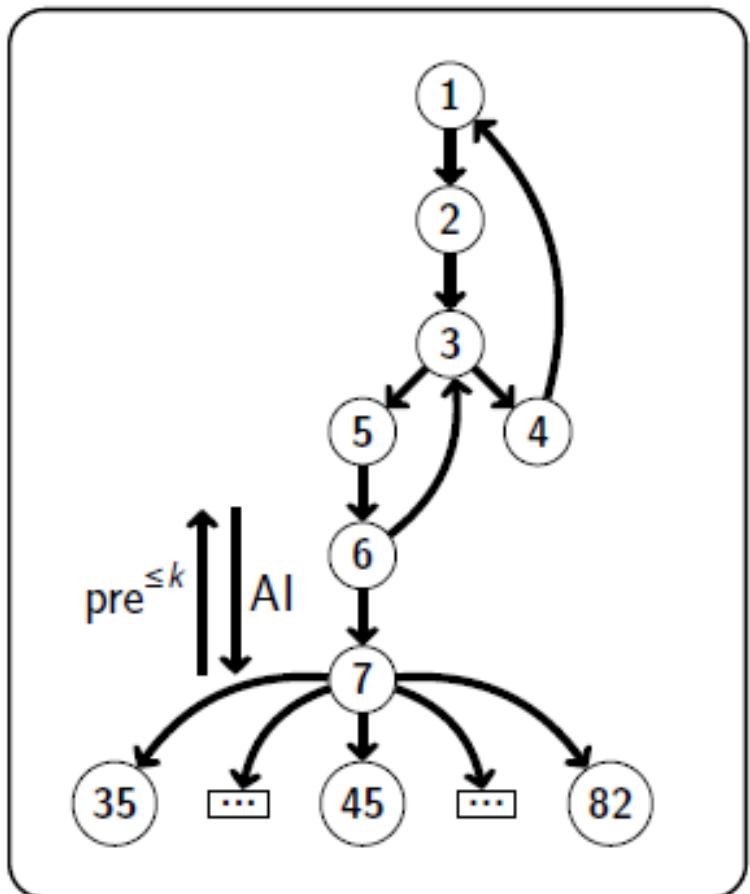
LOW-LEVEL CONDITIONS

	flag predicate	cmp x y predicate	sub x y predicate ²	test x y predicate
ja, jnbe	$\neg CF \wedge \neg ZF$	$x >_U y$	$x' \neq 0$	$x \& y \neq 0$
jae, jnb, jnc	$\neg CF$	$x \geq_U y$	true	true
jb, jnae, jc	CF	$x <_U y$	$x' \neq 0$	false
jbe, jna	$CF \vee ZF$	$x \leq_U y$	true	$x \& y = 0$
je, jz	ZF	$x = y$	$x' = 0$	$x \& y = 0$
jne, jnz	$\neg ZF$	$x \neq y$	$x' \neq 0$	$x \& y \neq 0$
jg, jnle	$\neg ZF \wedge (OF = SF)$	$x > y$	$x' > 0$	$(x \& y \neq 0) \wedge (x \geq 0 \vee y \geq 0)$
jge, jnl	$(OF = SF)$	$x \geq y$	true	$(x \geq 0 \vee y \geq 0)$
jl, jnge	$(OF \neq SF)$	$x < y$	$x' < 0$	$(x < 0 \wedge y < 0)$
jle, jng	$ZF \vee (OF \neq SF)$	$x \leq y$	true	$(x \& y = 0) \vee (x < 0 \wedge y < 0)$

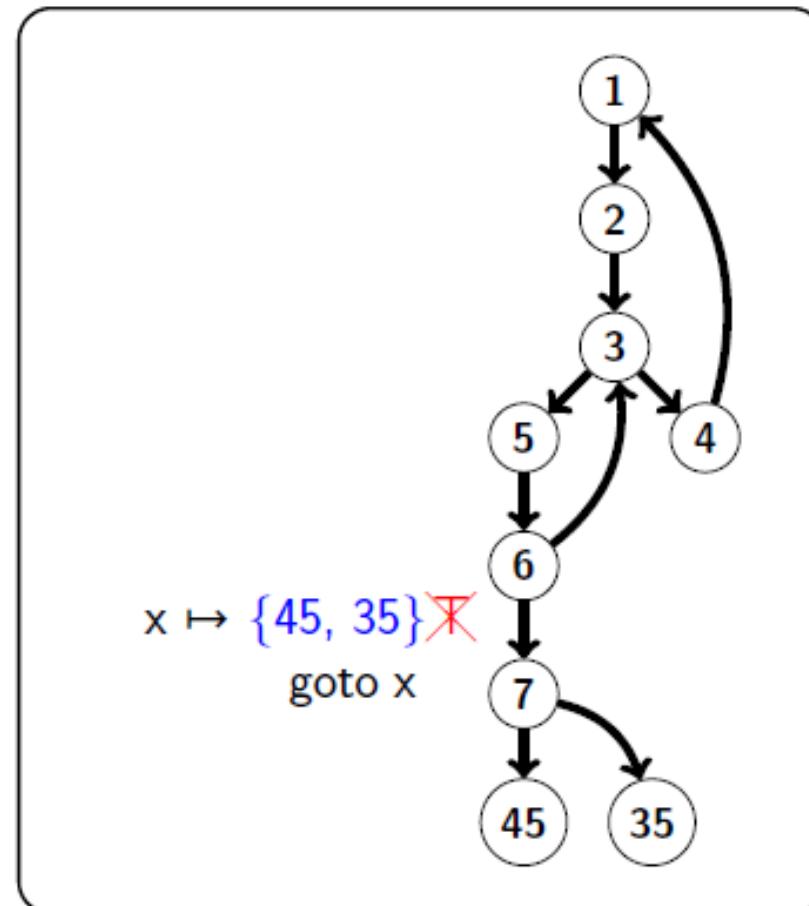
LOW-LEVEL CONDITIONS

example	retrieved condition	patterns
or eax, 0	if (<code>eax = 0</code>) then goto ...	✗
je ...		
cmp eax, 0	if (<code>eax ≥ 0</code>) then goto ...	✗
jns ...		
sar ebp, 1	if (<code>ebp = 0</code>) then goto ...	✗
je ...		
dec ecx	if (<code>ecx ≥ 0</code>) then goto ...	✗
jg ...		

Precision refinement [Brauer, 2011]



Degraded mode [Kinder, 2012]



Insights

- Complex predicates often hide simple predicates
- Only a few templates : $>_{u,s}$, $<_{u,s}$, $\geq_{u,s}$, $\leq_{u,s}$, $=$, \neq
- Try to find the appropriate one through equivalence checking
- Optimization :
 - Once per address using cache
 - Cheap pruning through filtering

Approach	archi. independent	Sound	Complete enough
Patterns	✗	✓ / ✗	✗
Logic-based	✓	✓	✗
Template-based	✓	✓	✓

HIGH-LEVEL CONDITION RECOVERY

method	#loc [†]	#cond [‡]	#success [*]	time	time _{all}
templates	242884	1978	1760 (89%)	22.93	2674.81
logic-based	247894	2260	694 (31%)	0.003	2561.64
patterns	229255	1987	1357 (68%)	0.014	2373.33
templates+patterns	242884	1978	1838 (92%)	9.17	2659.95

STATIC ANALYSIS in BINSEC

an overview

	difficulty	solution
Domains	low-level arithmetic ubiquitous data moves low-level conditions	dual-intervals equality domain flag domain + condition recovery
Widening	no loop structure	<i>loop detection</i> widening point positioning
CFG	unavailable recovery scale	incremental CFG recovery backward precision recovery degraded mode [Kinder2012]

	Correct	Complete	Efficient	Robust
Static syntactic	X	X / --	OK	X
Dynamic	OK	XX	OK	OK
DSE	OK	--	X	OK
Static semantic	X	OK / X	X	X

- Why binary-level analysis?
- Some background on source-level formal methods
- The hard journey from source to binary
- A few case-studies
- Conclusion

APPLICATION: VULNERABILITY DETECTION

Find vulnerabilities before the bad guys

- On the whole program
- At binary-level
- Know only the entry point and program input format



```
4800 6000 5dc3 5589 e5c7 0812 6000 00b8 4800 6000 5dc3 558  
0000 00b8 4500 0000 5dc3 5589 e5c7 0812 6000 00b8 4800 6000  
bf0e 0821 0000 00b8 5dc3 5589 e5c7 0812 6000 00b8 4500 000  
e5c7 0540 bf0e 0822 0000 00b8 5dc3 5589 e5c7 0540 bf0e 082  
5dc3 5589 e583 ec10 c705 00b8 4900 0000 5dc3 5589 e583 ec1  
0000 a148 bf0e 0883 f809 48bf 6e08 0100 0000 a148 bf0e 088  
8b04 8548 e10b 08FF e0c6 0597 0002 0000 8b04 8548 e10b 08F  
00c6 45f9 00c6 45fa 00c7 45f7 00c6 45f8 00c6 45f9 00c6 45f  
0000 00e9 d981 0000 c645 0548 bf0e 0882 0000 00e9 d981 000  
c645 f900 c645 fa01 807d f701 c645 f800 c645 f900 c645 fa0  
48bf 6e08 0300 0000 807d f800 750a c705 48bf 6e08 0300 000  
fc00 750a c705 48bf 6e08 f800 7410 807d fc00 750a c705 48b  
fc00 7415 807d fb00 740f 0900 6000 807d fc00 7415 807d fb0  
0600 0000 c988 0100 00e9 c705 48bf 6e08 0600 0000 c988 010  
f701 c645 f800 c645 f900 0000 0000 c645 f701 c645 f800 c64  
fc00 740f c705 48bf 6e08 c645 fa02 807d fc00 740f c705 48b  
0100 00e9 5901 0000 c645 0400 6000 e95e 0100 00e9 5901 000  
c645 f900 c645 fa03 807d f701 c645 f800 c645 f900 c645 fa0  
fe00 750a c705 48bf 6e08 f800 7410 807d fe00 750a c705 48b  
fc00 750a c705 48bf 6e08 0500 6000 807d fc00 750a c705 48b  
fe00 740f c705 48bf 6e08 0300 6000 807d fc00 740f c705 48b  
0100 0000 c645 0600 6000 e90e 0100 00e9 0901 000  
c645 0043 0001 807d f701 c645 f800 c645 f901 c645 fa0  
48bf 0400 0000 c9e4 0000 0000 750f c705 48bf 6e08 0400 000  
0000 c645 f701 c645 f800 0000 00e9 df00 0000 c645 f701 c64  
fa04 807d fc00 7410 807d c645 1900 c645 fa04 807d fc00 741  
48bf 6e08 0700 0000 807d ff00 750a c705 48bf 6e08 0700 000  
ff00 740f c705 48bf 6e08 fc00 7415 807d ff00 740f c705 48b  
0000 00e9 9900 0000 c645 0600 0000 e99e 0000 00e9 9900 000  
c645 f900 c645 fa05 807d f701 c645 f800 c645 f900 c645 fa0  
fc00 750a c705 48bf 6e08 f800 7410 807d fe00 750a c705 48b  
fc00 750a c705 48bf 6e08 0800 6000 807d fc00 750a c705 48b  
fe00 7506 807d ff00 740c 0900 6000 807d fe00 7506 807d ff0  
0600 6000 eb4b eb49 c645 c705 48bf 6e08 0600 6000 eb4b eb4  
c645 f901 c645 fa02 807d f701 c645 f800 c645 f901 c645 fa0  
5dc3 5589 e5c7 0540 bf0e 00b8 5400 0000 5dc3 5589 e5c7 054  
1800 6000 5dc3 5589 e5c7 0812 6000 00b8 4800 6000 5dc3 558  
3000 00b8 4500 0000 5dc3 0540 bf0e 0820 0000 0001  
bf0e 0821 0000 00b8 5800 5589 e5c7 0540 bf0e 082 USE 00b1  
e5c7 0540 bf0e 0822 0000 0000 5dc3 5589 e5c7 054 082  
5dc3 5589 e583 ec10 c705 00b8 4900 0000 5dc3 5589 e583 ec1  
3000 a148 bf0e 0883 f809 48bf 6e08 0100 0000 a148 bf0e 088  
3b04 8548 e10b 08FF e0c6 0F87 0002 0000 8b04 8548 e10b 08F  
00c6 45f9 00c6 45fa 00c7 45f7 00c6 45f8 00c6 45f9 00c6 45f  
0000 00e9 d981 0000 c645 0548 bf0e 0802 0000 00e9 d981 000  
c645 f900 c645 fa01 807d f701 c645 f800 c645 f900 c645 fa0  
48bf 6e08 0300 0000 807d fb00 750a c705 48bf 6e08 0300 000  
fc00 750a c705 48bf 6e08 f800 7410 807d fc00 750a c705 48b  
fc00 7415 807d fb00 740f 0900 6000 807d fc00 7415 807d fb0  
3000 0000 c988 0100 00e9 c705 48bf 6e08 0600 6000 c988 0101
```

APPLICATION: VULNERABILITY DETECTION

Many successful applications of pure DSE

- SAGE @ Microsoft
- Mayhem/VeriT @ ForallSecure

cf. Cyber Grand Challenge



mortderire.com

```
4800 6000 5dc3 5589 e5c7 0812 6000 00b8 4800 6000 5dc3 558  
0000 00b8 4500 0000 5dc3 5589 e5c7 0812 6000 00b8 4800 6000  
bf0e 0821 0000 00b8 4500 0000 5dc3 5589 e5c7 0812 6000 00b  
e5c7 0540 bf0e 0822 0000 00b8 4500 0000 5dc3 5589 e5c7 0540  
5dc3 5589 e583 ec10 c705 00b8 4900 0000 5dc3 5589 e583 ec1  
0000 a148 bf0e 0883 f809 48bf 6e08 0160 0000 a148 bf0e 088  
8b04 8548 e10b 08FF e0c6 0597 0002 0000 8b04 8548 e10b 08F  
00c6 45f9 00c6 45fa 00c7 45f7 00c6 45f8 00c6 45f9 00c6 45f  
0000 00e9 d981 0000 c645 0548 bf0e 0882 0000 00e9 d981 000  
c645 f900 c645 fa01 807d f701 c645 f800 c645 f900 c645 fa0  
48bf 6e08 0300 0000 807d fb00 750a c705 48bf 6e08 0300 000  
fc00 750a c705 48bf 6e08 fb00 7410 807d fc00 750a c705 48b  
fc00 7415 807d fb00 740f 0900 0000 807d fc00 7415 807d fb0  
0600 0000 c988 0100 00e9 c705 48bf 6e08 0600 0000 c988 010  
f701 c645 f800 c645 f900 0000 0000 c645 f701 c645 f800 c64  
fc00 740f c705 48bf 6e08 c645 fa02 807d fc00 740f c705 48b  
0100 00e9 5901 0000 c645 0400 0000 e95e 0100 00e9 5901 000  
c645 f900 c645 fa03 807d f701 c645 f800 c645 f900 c645 fa0  
fe00 750a c705 48bf 6e08 f000 750a c705 48bf 6e08 0000 000  
fc00 750a c705 48bf 6e08 0500 0000 807d fc00 750a c705 48b  
fe00 740f c705 48bf 6e08 0300 0000 807d fc00 740f c705 48b  
0100 0000 c645 0600 0000 e90e 0100 00e9 0901 000  
c645 0043 f001 807d f701 c645 f800 c645 f901 c645 fa0  
48bf 0400 0000 c9e4 f400 750f c705 48bf 6e08 0400 000  
0000 c645 f701 c645 f800 0000 00e9 df00 0000 c645 f701 c64  
fa04 807d fc00 7410 807d c645 1900 c645 fa04 807d fc00 741  
48bf 6e08 0700 0000 807d ff00 750a c705 48bf 6e08 0700 000  
ff00 740f c705 48bf 6e08 fc00 7415 807d ff00 740f c705 48b  
0000 00e9 9900 0000 c645 0600 0000 e99e 0000 00e9 9900 000  
c645 f900 c645 fa05 807d f701 c645 f800 c645 f900 c645 fa0  
fc00 750a c705 48bf 6e08 f000 7410 807d fe00 750a c705 48b  
fc00 750a c705 48bf 6e08 0800 0000 807d fc00 750a c705 48b  
fe00 7506 807d ff00 740c 0900 0000 807d fe00 7506 807d ff0  
0600 6000 eb4b eb49 c645 c705 48bf 6e08 0600 0000 eb4b eb4  
c645 f901 c645 fa02 807d f701 c645 f800 c645 f901 c645 fa0  
5dc3 5589 e5c7 0540 bf0e 00b8 5400 0000 5dc3 5589 e5c7 054  
1800 6000 5dc3 5589 e5c7 0812 6000 00b8 4800 6000 5dc3 558  
3000 00b8 4500 0000 5dc3 0540 bf0e 0820 0000 0000 0000 0000  
bf0e 0821 0000 00b8 5800 5589 e5c7 0540 bf0e 0821 0000 0000  
e5c7 0540 bf0e 0822 0000 0000 5dc3 5589 e5c7 0540 bf0e 0822  
5dc3 5589 e583 ec10 c705 00b8 4900 0000 5dc3 5589 e583 ec1  
3000 a148 bf0e 0883 f809 48bf 6e08 0100 0000 a148 bf0e 088  
3b04 8548 e10b 08FF e0c6 0F87 0002 0000 8b04 8548 e10b 08F  
00c6 45f9 00c6 45fa 00c7 45f7 00c6 45f8 00c6 45f9 00c6 45f  
0000 00e9 d981 0000 c645 0548 bf0e 0882 0000 00e9 d981 0000  
c645 f900 c645 fa01 807d f701 c645 f800 c645 f900 c645 fa0  
48bf 6e08 0300 0000 807d fb00 750a c705 48bf 6e08 0300 0000  
fc00 750a c705 48bf 6e08 fb00 7410 807d fc00 750a c705 48b  
fc00 7415 807d fb00 740f 0900 0000 807d fc00 7415 807d fb0  
3000 0000 c988 0100 00e9 c705 48bf 6e08 0600 0000 c988 0100
```

APPLICATION: VULNERABILITY DETECTION

[SSPREW 2016, with VERIMAG]

Here:

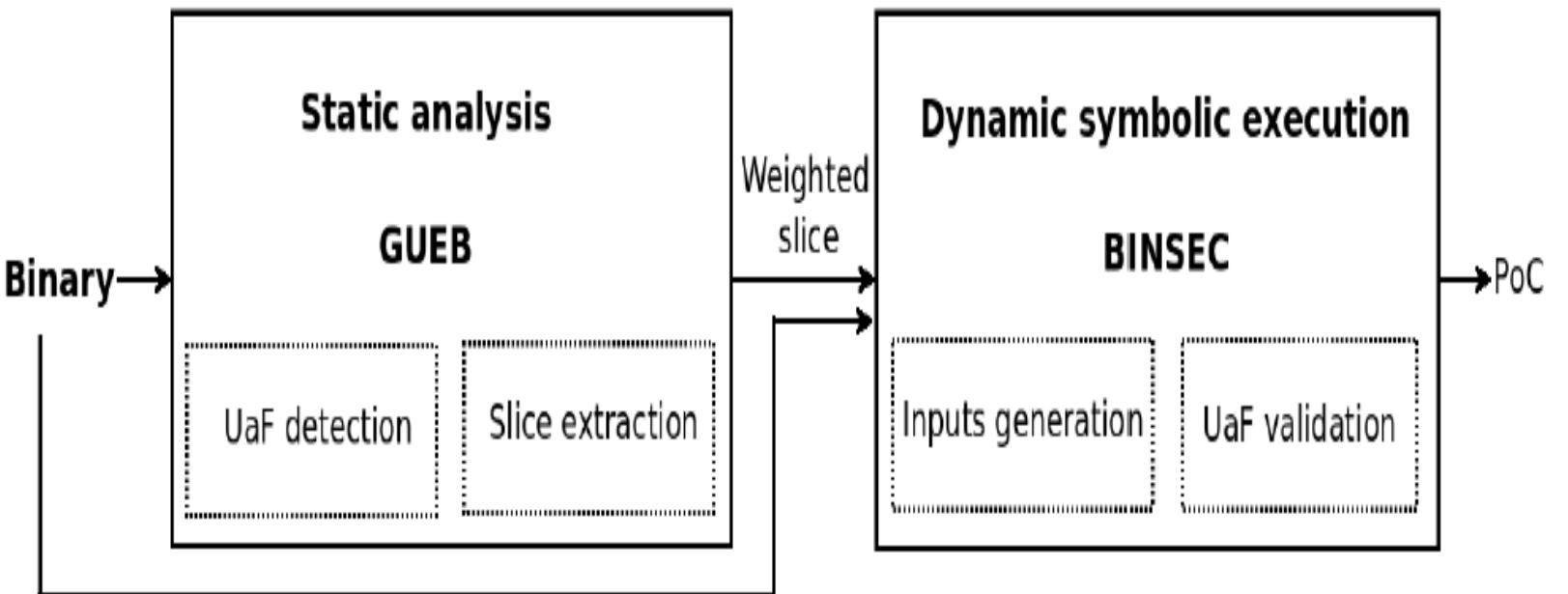
- Focus on use-after-free
- Combine static and DSE



mortdeirine.com

```
4800 6000 5dc3 5589 e5c7 0812 6000 00b8 4800 6000 5dc3 558  
0000 00b8 4500 0000 5dc3 5589 e5c7 0812 6000 00b8 4800 6000  
bf0e 0821 0000 00b8 5dc3 5589 e5c7 0812 6000 00b8 4500 000  
e5c7 0540 bf0e 0822 0000 00b8 5dc3 5589 e5c7 0540 bf0e 082  
5dc3 5589 e583 ec10 c705 00b8 4900 0000 5dc3 5589 e583 ec1  
0000 a148 bf0e 0883 f809 48bf 6e08 0100 0000 a148 bf0e 088  
8b04 8548 e10b 08FF e0c6 0597 0002 0000 8b04 8548 e10b 08F  
00c6 45f9 00c6 45fa 00c7 45f7 00c6 45f8 00c6 45f9 00c6 45f  
0000 00e9 d981 0000 c645 0548 bf0e 0882 0000 00e9 d981 000  
c645 f900 c645 fa01 807d f701 c645 f800 c645 f900 c645 fa0  
48bf 6e08 0300 0000 807d fb00 750a c705 48bf 6e08 0300 000  
fc00 750a c705 48bf 6e08 fb00 7410 807d fc00 750a c705 48b  
fc00 7415 807d fb00 740f 0900 0000 807d fc00 7415 807d fb0  
0600 0000 c988 0100 00e9 c705 48bf 6e08 0600 0000 c988 010  
f701 c645 f800 c645 f900 0000 0000 c645 f701 c645 f800 c64  
fc00 740f c705 48bf 6e08 c645 fa02 807d fc00 740f c705 48b  
0100 00e9 5901 0000 c645 0400 0000 e95e 0100 00e9 5901 000  
c645 f900 c645 fa03 807d f701 c645 f800 c645 f900 c645 fa0  
fe00 750a c705 48bf 6e08 f000 750a c705 48bf 6e08 0100 00e9 5901 000  
fc00 750a c705 48bf 6e08 0500 0000 807d fc00 750a c705 48b  
fe00 740f c705 48bf 6e08 0300 0000 807d fc00 740f c705 48b  
0100 0000 c645 0600 0000 e90e 0100 00e9 0901 000  
c645 0043 f001 807d f701 c645 f800 c645 f901 c645 fa0  
48bf 0400 0000 c9e4 f400 750f c705 48bf 6e08 0400 000  
0000 c645 f701 c645 f800 0000 00e9 df00 0000 c645 f701 c64  
fa04 807d fc00 7410 807d c645 f900 c645 fa04 807d fc00 741  
48bf 6e08 0700 0000 807d ff00 750a c705 48bf 6e08 0700 000  
ff00 740f c705 48bf 6e08 fc00 7415 807d ff00 740f c705 48b  
0000 00e9 9900 0000 c645 0600 0000 e99e 0000 00e9 9900 000  
c645 f900 c645 fa05 807d f701 c645 f800 c645 f900 c645 fa0  
fc00 7500 c705 48bf 6e08 f000 7410 807d fe00 750a c705 48b  
fc00 750a c705 48bf 6e08 0800 0000 807d fc00 750a c705 48b  
fe00 7506 807d ff00 740c 0900 0000 807d fe00 7506 807d ff0  
0600 6000 eb4b eb49 c645 c705 48bf 6e08 0600 0000 eb4b eb4  
c645 f901 c645 fa02 807d f701 c645 f800 c645 f901 c645 fa0  
5dc3 5589 e5c7 0540 bf0e 00b8 5400 0000 5dc3 5589 e5c7 054  
1800 6000 5dc3 5589 e5c7 0812 6000 00b8 4800 6000 5dc3 558  
3000 00b8 4500 0000 5dc3 0540 bf0e 0820 0000 0001  
bf0e 0821 0000 00b8 5800 5589 e5c7 0540 bf0e 082 USE 00b1  
e5c7 0540 bf0e 0822 0000 0000 5dc3 5589 e5c7 054 0821  
5dc3 5589 e583 ec10 c705 00b8 4900 0000 5dc3 5589 e583 ec11  
3000 a148 bf0e 0883 f809 48bf 6e08 0100 0000 a148 bf0e 088  
3b04 8548 e10b 08FF e0c6 0F87 0002 0000 8b04 8548 e10b 08F  
00c6 45f9 00c6 45fa 00c7 45f7 00c6 45f8 00c6 45f9 00c6 45f  
0000 00e9 d981 0000 c645 0548 bf0e 0882 0000 00e9 d981 000  
c645 f900 c645 fa01 807d f701 c645 f800 c645 f900 c645 fa0  
48bf 6e08 0300 0000 807d fb00 750a c705 48bf 6e08 0300 000  
fc00 750a c705 48bf 6e08 fb00 7410 807d fc00 750a c705 48b  
fc00 7415 807d fb00 740f 0900 0000 807d fc00 7415 807d fb0  
3000 0000 c988 0100 00e9 c705 48bf 6e08 0600 0000 c988 0101
```

KEY IDEAS (Josselin Feist)



A Pragmatic 2-step approach

- Static:** scale, not complete, not correct
- Symbolic:** correct, directed by static
- Combination:** scalable and correct

```

4800 0000 5dc3 5589 e5c7 0812 0000 00b8 4800 0000 5dc3 558
0000 00b8 4500 0000 5dc3 5589 e5c7 0812 0000 00b8 4800 0000 5dc3 558
bf0e 0821 0000 00b8 5dc3 5589 e5c7 0812 0000 00b8 4800 0000 5dc3 558
e5c7 0540 bf0e 0822 0000 00b8 5dc3 5589 e5c7 0812 0000 00b8 4800 0000 5dc3 558
5dc3 5589 e583 ec10 c705 00b8 4900 0000 5dc3 5589 e583 ec1
0000 a148 bf0e 0883 f809 48bf 0e08 0100 0000 a148 bf0e 088
8b04 8548 e10b 08ff e0c6 0f87 0002 0000 8b04 8548 e10b 08f
00c6 45f9 00c6 45fa 00c7 45f7 00c6 45fb 00c6 45f9 00c6 45f
0000 00e9 d901 0000 c645 0548 bf0e 0000 0000 00e9 d901 000
c645 f900 c645 fa01 807d 7701 c645 f800 c645 f900 c645 fa0
48bf 0e68 0300 0000 807d fb00 750a c705 48bf 0e68 0300 000
fc00 750a c705 48bf 0e08 fb00 7410 807d fc00 750a c705 48b
fc00 7415 807d fb00 740f 0900 0000 807d fc00 7415 807d fb0
0600 0000 c988 0100 00e9 c705 48f1 0e08 0600 0000 c988 010
f701 c645 f800 c645 f900 8301 0000 c645 f701 c645 f800
fc00 745f c705 48f2 0e08 c645 fa02 807d fc00 740f c705 48b
0100 00e9 5901 0000 c645 0400 0000 e95e 0100 00e9 5901 000
c645 f900 c645 fa03 807d f701 c645 f800 c645 f900 c645 fa0
fe00 7502 c705 48bf 0e08 fd00 7410 807d fe00 750a c705 48b
fc00 7502 c705 48bf 0e08 0500 0000 807d f700 750a c705 48b
fe00 740f c705 48bf 0e08 0300 0000 807d fe00 740f c705 48b
0100 1901 0000 c645 0600 0000 e902 0100 00e9 0901 000
c645 0445 f701 807d f701 c645 1800 c645 f901 c645 fa0
48bf 0400 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000
0000 c645 f701 c645 f800 0000 00e9 d100 0000 c645 f701 c64
5004 807d fc00 7410 807d c645 f800 c645 fa04 807d fc00 741
48bf 0e68 0700 0000 807d ff00 750a c705 48bf 0e68 0700 000
ff00 740f c705 48bf 0e08 c645 f900 7415 807d fc00 740f c705 48b
0000 00e9 9900 0000 c645 0600 0000 e99e 0000 00e9 9900 000
c645 f900 c645 fa05 807d f701 c645 f800 c645 f900 c645 fa0
fe00 750a c705 48bf 0e08 f100 7410 807d fe00 750a c705 48b
fc00 750a c705 48bf 0e08 f100 7410 807d fe00 750a c705 48b
fe00 7506 807d ff00 740c 0800 0000 807d fc00 750a c705 48b
0600 0000 eb4b eb49 c645 c705 48bf 0e08 0600 0000 eb4b eb4
c645 f901 c645 fa02 807d f701 c645 f800 c645 f901 c645 fa0
5dc3 5589 e5c7 0540 bf0e 00b8 5400 0000 5dc3 5589 e587 0541
4800 0000 5dc3 5589 e5c7 0812 0000 00b8 4800 0000 5dc3 5581
3000 00b8 4500 0000 5dc3 0540 bf0e 0100 0000 0000 0000 0001
bf0e 0821 0000 00b8 5800 5589 e5c7 0541 bf0e 0821 0000 0001
e5c7 0540 bf0e 0822 0000 0000 5dc3 5589 e5c7 0541 bf0e 0821
5dc3 5589 e583 ec10 c705 00b8 4900 0000 5dc3 5589 e583 ec1
3000 a148 bf0e 0883 f809 48bf 0e08 0100 0000 a148 bf0e 088
8b04 8548 e10b 08ff e0c6 0f87 0002 0000 8b04 8548 e10b 08f
00c6 45f9 00c6 45fa 00c7 45f7 00c6 45fb 00c6 45f9 00c6 45f
0000 00e9 d901 0000 c645 0548 bf0e 0000 0000 00e9 d901 000
c645 f900 c645 fa01 807d f701 c645 f800 c645 f900 c645 fa0
48bf 0e68 0300 0000 807d fb00 750a c705 48bf 0e68 0300 000
fc00 750a c705 48bf 0e08 fb00 7410 807d fc00 750a c705 48b
fc00 7415 807d fb00 740f 0900 0000 807d fc00 7415 807d fb0
3600 0000 e988 0100 00e9 c705 48bf 0e08 0600 0000 e988 010

```

EXPERIMENTAL EVALUATION

■ GUEB + manual analysis [j. comp. virology 14]

- ▶ tiff2pdf : CVE-2013-4232
- ▶ openjpeg : CVE-2015-8871
- ▶ gifcolor : CVE-2016-3177
- ▶ accel-ppp

■ GUEB + BINSE/SE [ssprew16]

- ▶ Jasper JPEG-2000 : CVE-2015-5221

On these examples:

- Better than DSE alone
- Better than blackbox fuzzing
- Better than greybox fuzzing with no seed

```
48000 00000 5dc3 5589 e5c7 / 0812 0000 00b8 4800 0000 5dc3 558
0000 00b8 4500 0000 5 1820 0000 00b8 4500 000
bf0e 0821 0000 00b8 4500 0000 5dc3 5589 e5c7 0540 000
e5c7 0540 bf0e 0822 0000 00b8 4500 0000 5dc3 5589 e5c7 0540 bf0e 082
5dc3 5589 e583 ec10 C70> 00b8 4900 0000 5dc3 5589 e583 ec1
0000 a148 bf0e 0883 f809 48bf 0e08 0100 0000 a148 bf0e 088
8b04 8548 e10b 08ff e0c6 0f87 0002 0000 8b04 8548 e10b 08f
00c6 45f9 00c6 45fa 00c7 45f7 00c6 45fb 00c6 45f
0000 00e9 d901 0000 c645 0548 bf0e 0002 0000 00e9 d901 000
c645 f900 c645 fa01 807d 1701 c645 f800 c645 f900 c645 fa0
48bf 0e08 0300 0000 807d fb00 750a c705 48bf 0e08 0300 000
fc00 750a c705 48bf 0e08 fb00 7410 807d fc00 750a c705 48b
fc00 7415 807d fb00 740f 0900 0000 807d fc00 7414 807d fb0
0600 0000 c988 0100 00e9 c705 49f 0e08 0600 0000 c988 010
f701 c645 f800 c645 f900 8301 0000 c645 f701 c645 f800 c64
fc00 740f c705 48bf 0e08 c645 fa02 807d fc00 740f c705 48b
0100 00e9 5901 0000 c645 0400 0000 e95e 0100 00e9 5901 000
c645 f900 c645 fa03 807d f701 c645 f800 c645 f900 c645 fa0
fe00 750a c705 48bf 0e08 fd00 7410 807d fe00 750a c705 48b
fc00 750a c705 48bf 0e08 0500 0000 807d fe00 750a c705 48b
fe00 740f c705 48bf 0e08 0300 0000 807d fe00 740f c705 48b
0100 0901 0000 c645 0600 0000 e90e 0100 00e9 0901 000
c645 045 f001 807d f701 343 f800 c645 f901 c645 fa0
48bf 0e08 0400 0000 807d fe00 740f c705 48bf 0e08 0400 000
0000 c645 f701 c645 f800 0600 0000 807d fe00 750a c705 48b
fc00 807d fc00 7410 807d c645 f900 c645 fa04 807d fc00 741
48bf 0e08 0700 0000 807d ff00 750a c705 48bf 0e08 0700 000
ff00 740f c705 48bf 0e08 fc00 7415 807d fe00 740f c705 48b
0000 00e9 9900 0000 c645 0600 0000 e99e 0000 00e9 9900 000
c645 f900 c645 fa05 807d f701 c645 f800 c645 f900 c645 fa0
fe00 750a c705 48bf 0e08 f100 7410 807d fe00 750a c705 48b
0e08 0800 0000 807d fc00 740a c705 48b
740c 0900 0000 807d fe00 750a 807d ff00
c645 c705 48bf 0e08 0600 0000 cb4b eb4
807d f701 c645 f800 c645 f901 c645 fa0
bf0e 00b8 5400 0000 5dc3 5589 e583 0541
e5c7 0812 0000 00b8 4800 0000 5dc3 5589
5dc3 0540 bf0e 0820 0000 00b8 4800 0000
5800 5589 e5c7 0540 bf0e 0820 0000 00b8 4800 0000
0000 0000 5dc3 5589 e5c7 0540 bf0e 0820 0000 00b8 4800 0000
c705 00b8 4900 0000 5dc3 5589 e583 ec10
f809 48bf 0e08 0100 0000 a148 bf0e 0881
e0c6 0f87 0002 0000 8b04 8548 e10b 08f
00c7 45f7 00c6 45f8 00c6 15f9 00c6 45f
c645 0548 bf0e 0882 0000 00c9 d901 000
807d f701 c645 f800 c645 f900 c645 fa0
807d fb00 750a c705 48bf 0e08 0300 0000
0e08 fb00 7410 807d fc00 750a c705 48b
740f 0900 0000 807d fe00 7415 807d fb0
00e9 c705 48bf 0e08 0600 0000 c988 0100
```



APPLICATION: MALWARE DEOBFUSCATION

[S&P 2017, with LORIA]

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

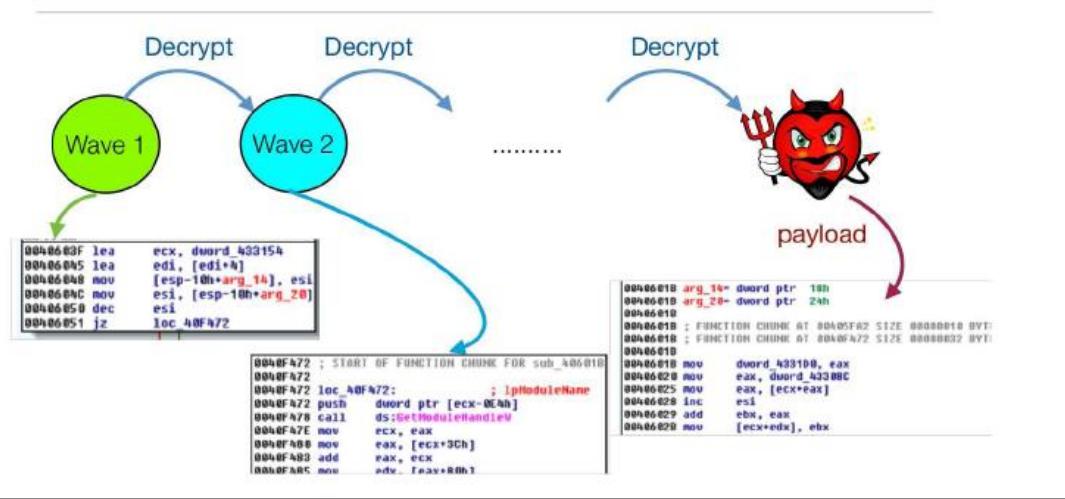
- understand what has been going on
- mitigate, fix and clean
- improve defense



Goal: help malware comprehension

- Reverse of heavily obfuscated code
- Identify and simplify protections

REVERSE CAN BECOME A NIGHTMARE (OBFUSCATION)



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

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

```

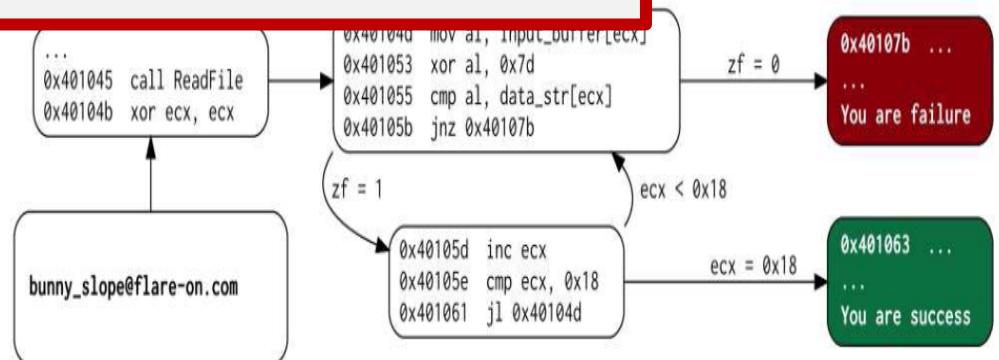
    mov eax, ds:X
    mov ecx, ds:Y
    imul ecx, ecx
    imul ecx, 7
    sub ecx, 1
    imul eax, eax
    cmp ecx, eax
    jz <dead_addr>
  
```

- Obfuscation**
hard to reverse
- self-modification
 - encryption
 - virtualization
 - code overlapping
 - opaque predicates
 - callstack tampering
 - ...

Goal: help malware comprehension

- Identify and simplify protections
- Ideal = revert protections

address	instr
80483d1	call +5
80483d6	pop edx
80483d7	add edx, 8
80483da	push edx
80483db	ret
80483dc	byte[invalid]



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)



```
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>
```

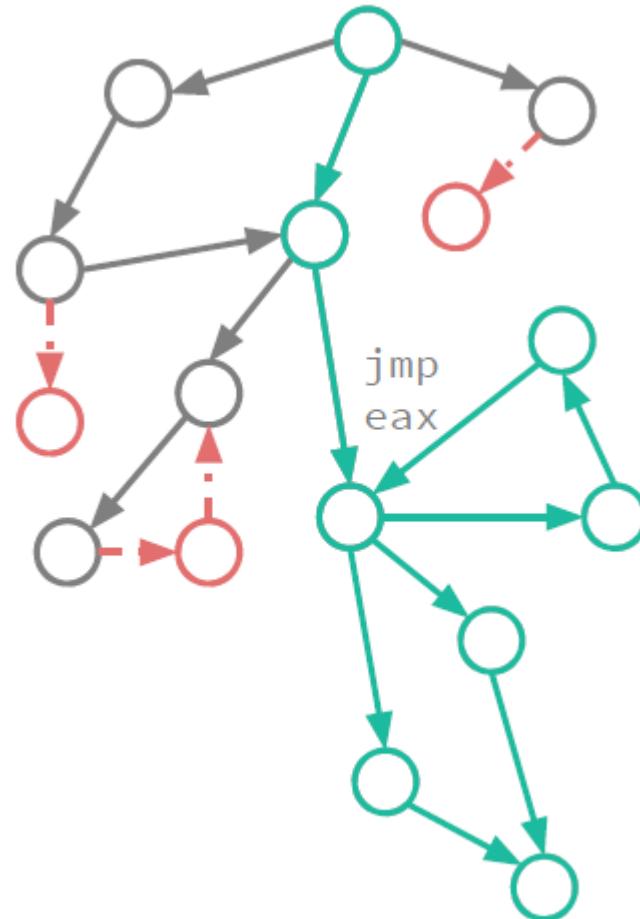
EXAMPLE: STACK TAMPERING

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

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

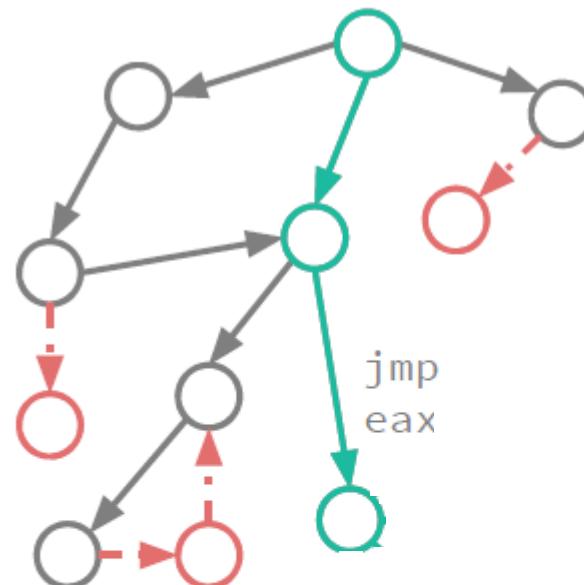
address	instr
80483d1	call +5
80483d6	pop edx
80483d7	add edx, 8
80483da	push edx
80483db	ret
80483dc	.byte{invalid}
80483de	[...]

STANDARD DISASSEMBLY TECHNIQUES ARE NOT ENOUGH



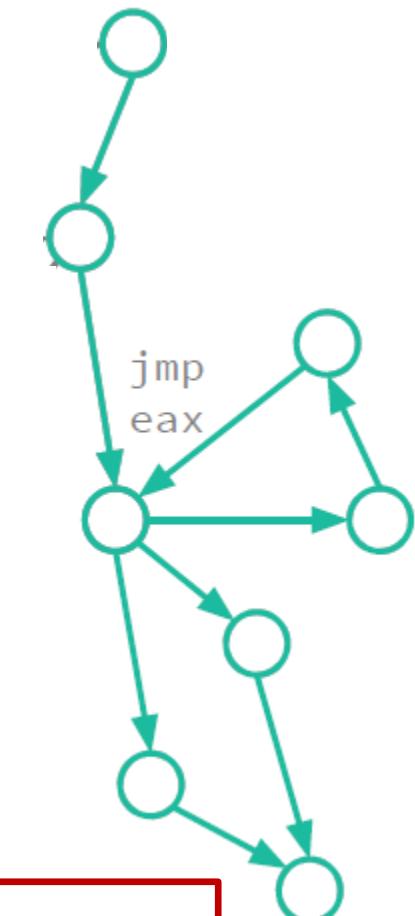
Static analysis

- too fragile vs obfuscation
- junk instr, missed instr.



Dynamic analysis

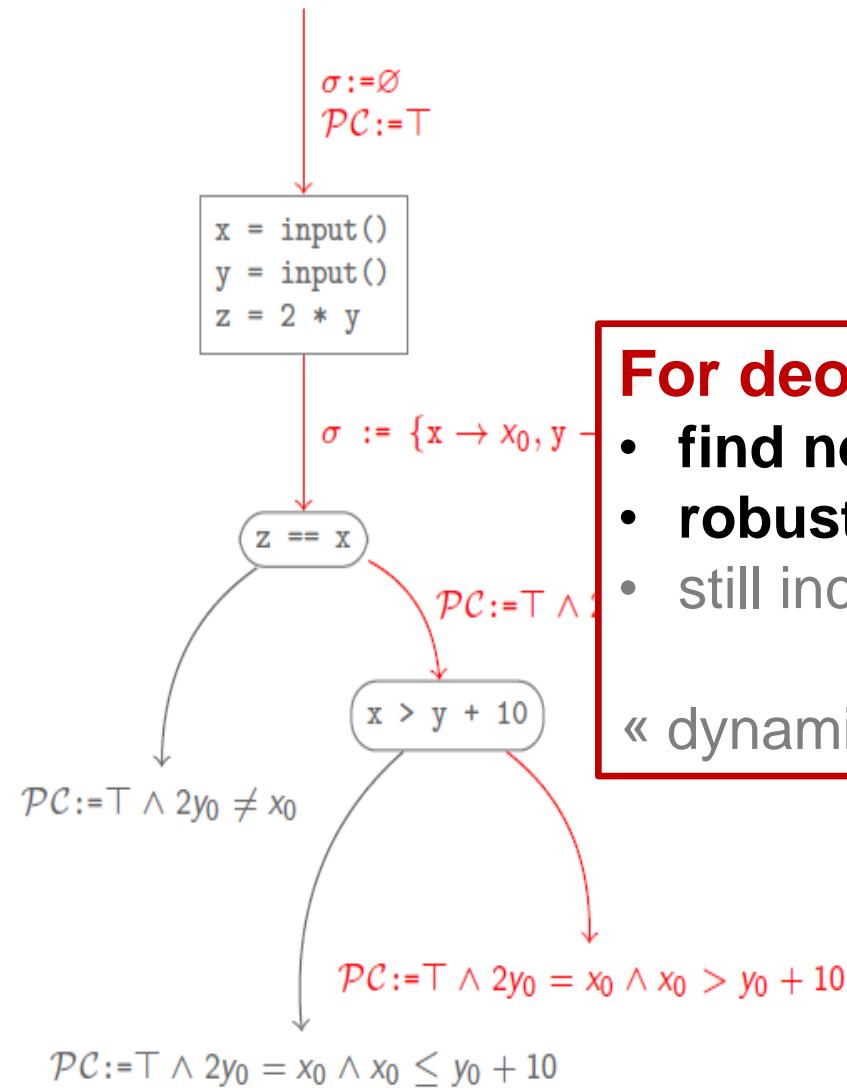
- robust vs obfuscation
- too incomplete



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 »

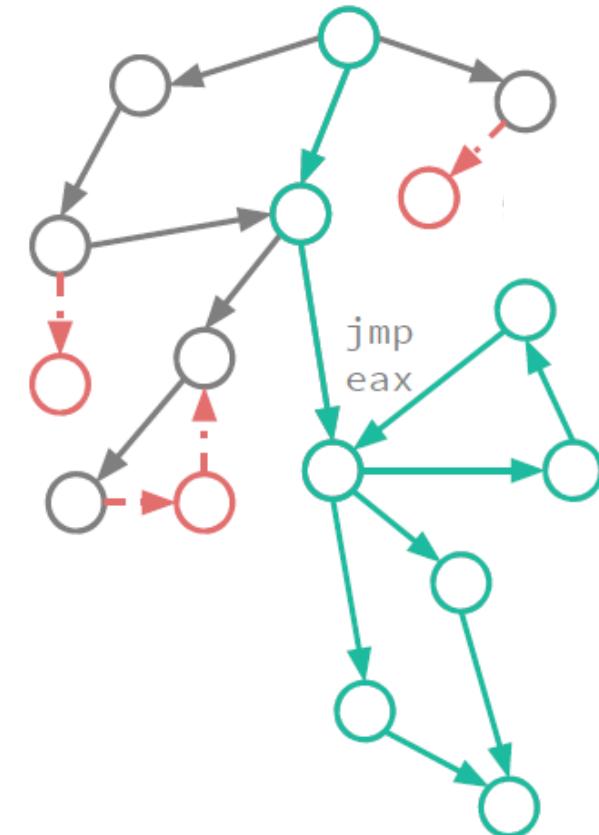
Prove that something is always true (resp. false)

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.?
 - ...



Not addressed by DSE

- Cannot enumerate all paths

OUR PROPOSAL: BACKWARD-BOUNDED SYMBOLIC EXECUTION

Insight 1: symbolic reasoning

- precision
- But: need finite #paths

Insight 2: backward-bounded

- $\text{pre}_k(c)=0 \Rightarrow c$ is infeasible
- finite #paths
- efficient, depends on k
- But: backward on jump eax?

Insight 3: dynamic partial CFG

- solve (partially) dyn. jumps
- robustness

Low FP/FN rates in practice

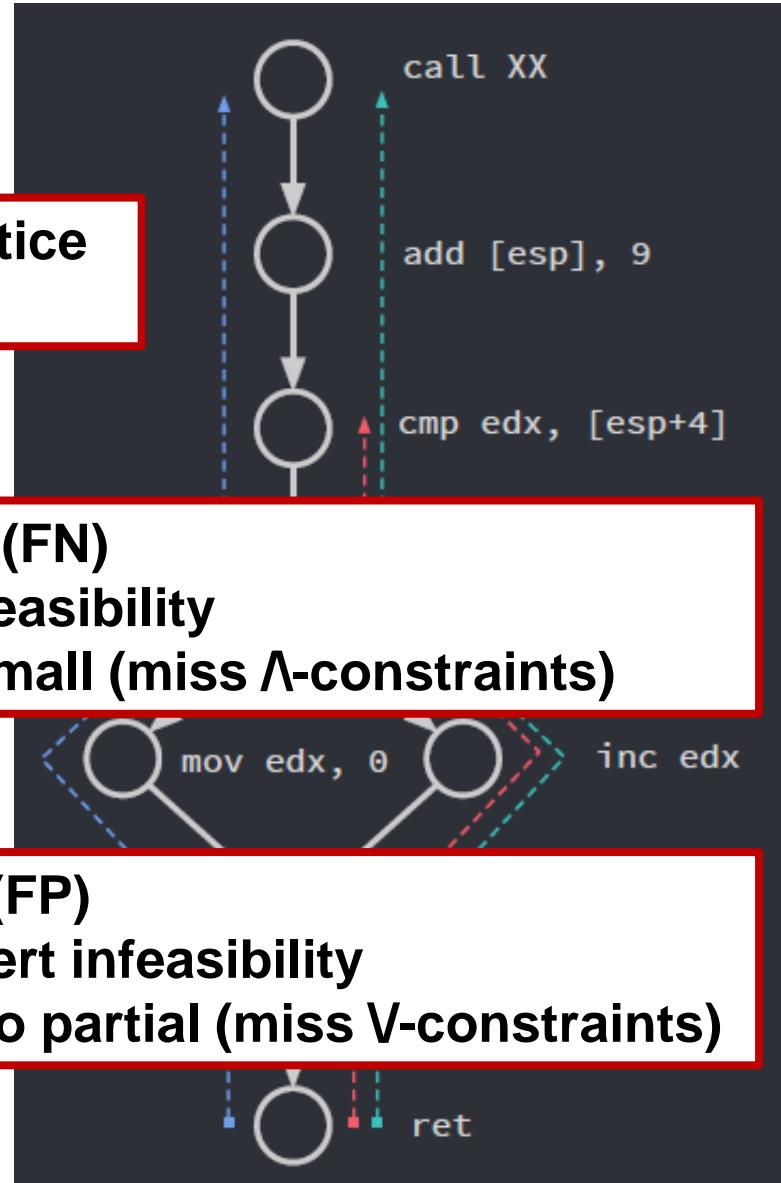
- ground truth xp

False negative (FN)

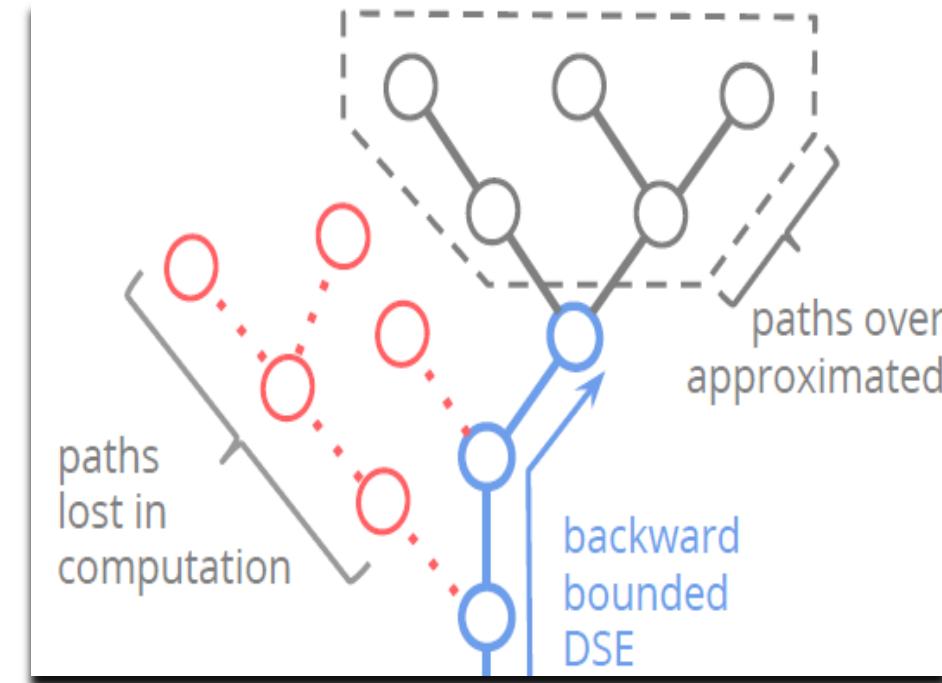
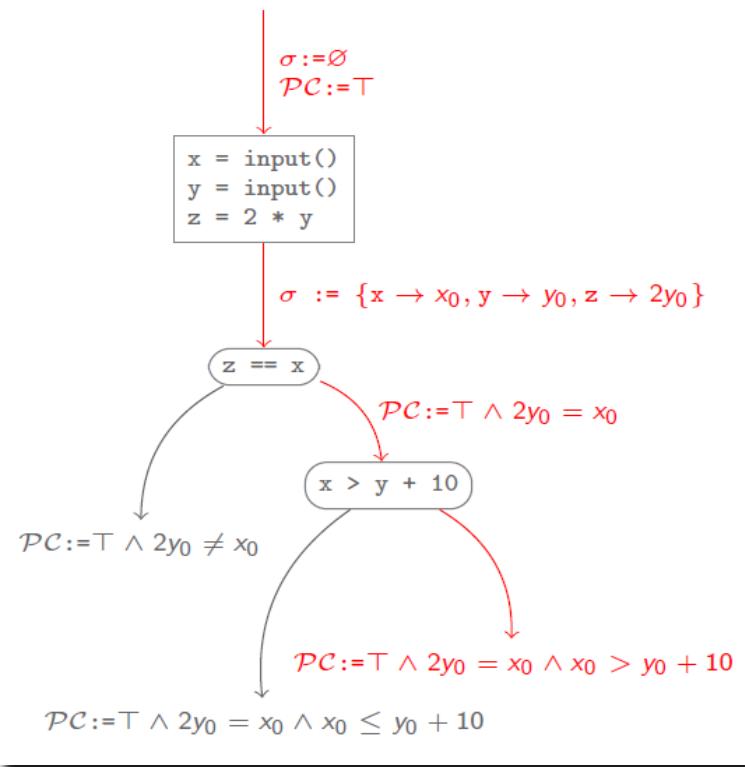
- can miss infeasibility
- why: k too small (miss \wedge -constraints)

False positive (FP)

- wrongly assert infeasibility
- why: CFG too partial (miss \vee -constraints)



FORWARD & BACKWARD SYMBOLIC EXECUTION



	(forward) DSE	bb-DSE
feasibility queries	●	●
infeasibility queries	●	●
scale	●	●

EXPERIMENTAL EVALUATION

- **Controlled experiments** (ground truth) → **precision**
- **Large-scale experiment: packers** → **scalability, robustness**
- **Case-study: X-tunnel malware** → **usefulness**

CONTROLLED EXPERIMENTS

- 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)		Genuine (5183)		TO	Error rate (FP+FN)/Tot (%)	Time (s)	avg/query (s)
	ok	miss (FN)	ok	miss (FP)				
2	0	5556	5182	1	0	51.75	89	0.008
4	292	4652	5152	20	0	42.61	96	0.009
8	102	4554	5152	20	0	39.89	120	0.011
16	102	4554	5152	20	0	38.83	152	0.014
32	5552	4	4579	604	25	5.66	197	0.018
40	5548	8	4523	660	39	6.22	272	0.025
50	5544	12	4458	725	79	6.86	384	0.036
							699	0.065
							1145	0.107
							2025	0.189

- Very precise résultats
- Seems efficient

Sample	runtime genuine			runtime violation		
	#ret †	proved genuine	proved a/d	#ret †	proved a/d	proved single
<i>obfuscated programs</i>						
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
<i>non-obfuscated programs</i>						
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

Obsidium
 JD Pack
 WinUpack PE Lock
 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

packers	trace len.	#proc	#th	#SMC	opaque predicates	call stack tampering
					OK	tamper
ACProtect v2.0	1.8M	1	1	1	159	0
ASPack v2.12	377K	1	1	1	24	11
Crypter v1.12	1.1M	1	1	1	24	125
Expressor	635K	1	1	1	14	0
FSG v2.0	68k	1	1	1	6	0
Mew	59K	1	1	1	6	0
PE Lock	2.3M	1	1	6	90	3
RLPack	941K	1	1	1	14	0
TELock v0.51	406K	1	1	1	3	1
Upack v0.39	711K	1	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
WinUpack
PE Lock
Expressor
PE Compact
Armadillo
Packman
EP Protector
ACProtect
TELockSVK
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
10043b0 jnz 0x10043ba
```

Enter SMC Layer 1

```
10043ba push 0x10011d7
10043bf ret
```

OP (decoy) in ASPack

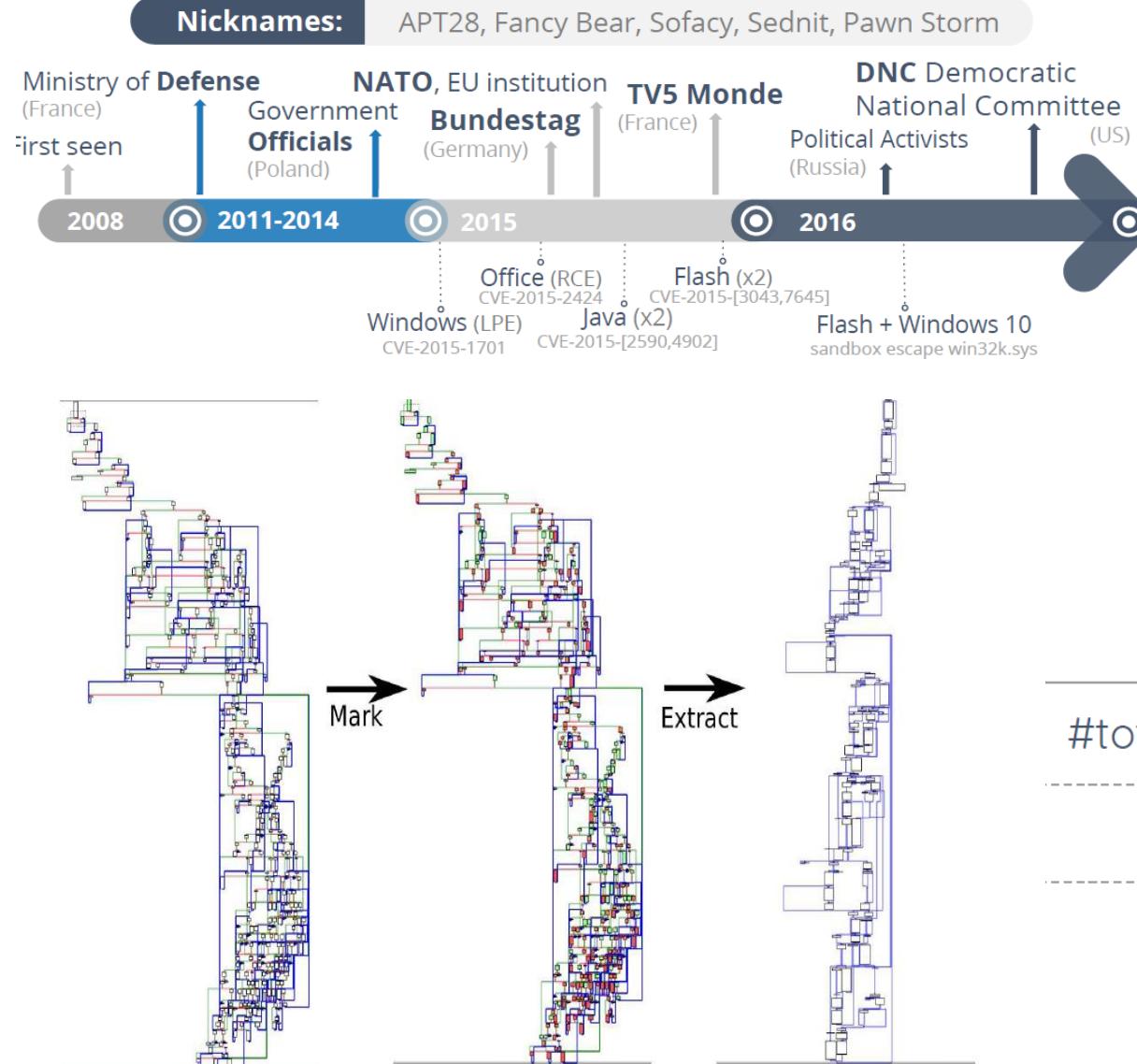
```
10040fe: mov bl, 0x0
10041c0: cmp bl, 0x1
1004103: jnz 0x1004163
```

ZF = 0 ZF = 1

```
1004163: jmp 0x100416d
1004105: inc [ebp+0xec]
[...]
```

0x10040ff at runtime

CASE-STUDY: THE XTUNNEL MALWARE (part of DNC hack)



Two heavily obfuscated samples

- Many opaque predicates

Goal: detect & remove protections

- Identify 50% of code as spurious
- Fully automatic, < 3h

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

$$7y^2 - 1 \neq x^2 \quad \frac{2}{x^2 + 1} \neq y^2 + 3$$

SECURITY ANALYSIS: COUNTER-MEASURES (and mitigations)

- **Long dependency 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**

Also

- « **Probabilistic obfuscation** »
- **Covert channels**

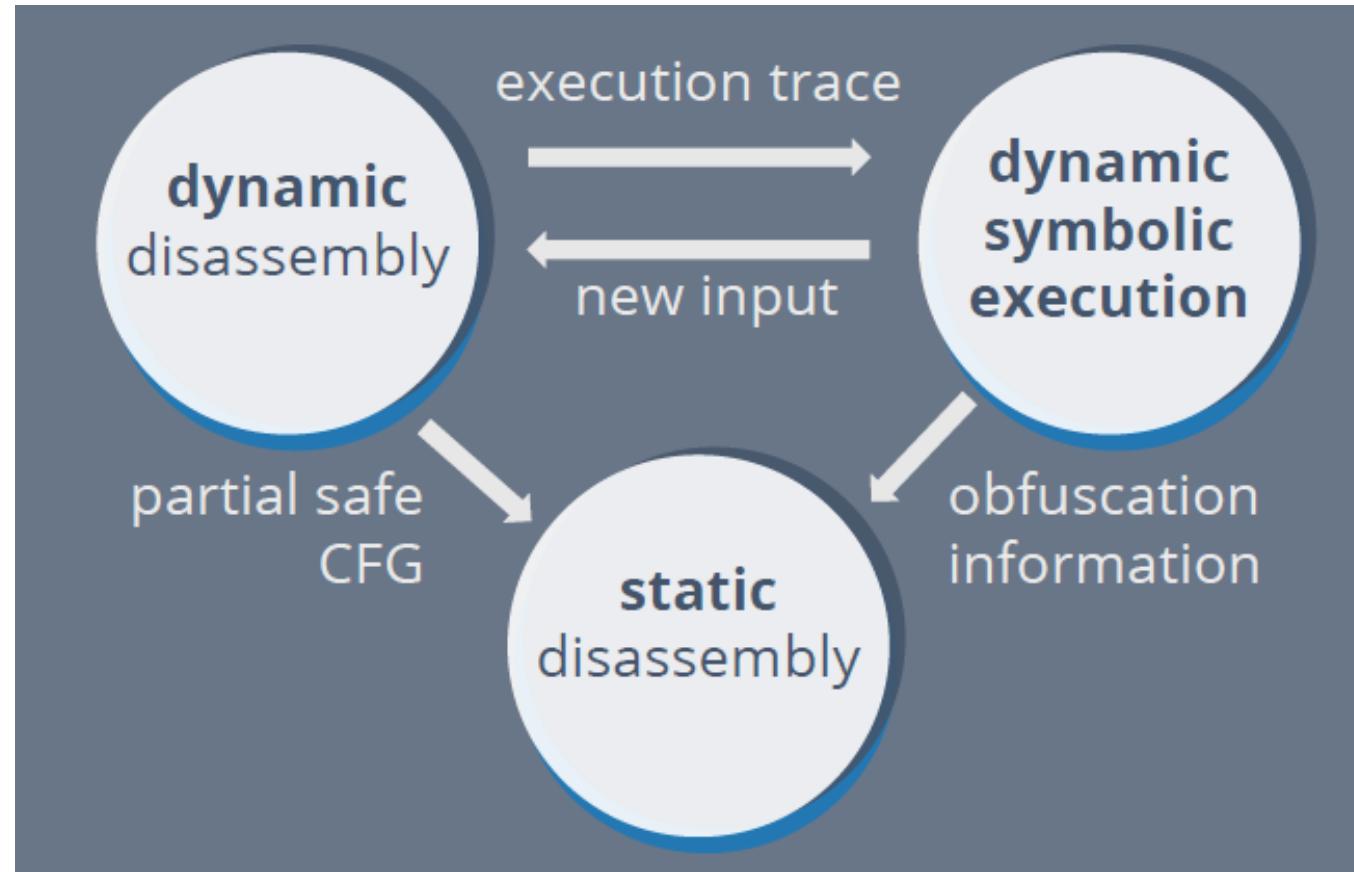
- Why binary-level analysis?
- Some background on source-level formal methods
- The hard journey from source to binary
- A few case-studies
- Conclusion

	Feasibility	Infeasibility	Efficient	Robust
Static syntactic	X	X	OK	X
Dynamic	--	X	OK	OK
DSE	OK	X	X	OK
Static semantic	X	OK	X	X
BB-DSE	X	OK (fp,fn)	OK	OK

- **Semantic analysis can change the game of binary-level security**
 - Current syntactic and dynamic methods are not enough
 - [complement existing approaches and help the expert, not replace everything]
 - Explore more, Prove invariance, Simplify
- **Yet, challenging to adapt from source-level safety-critical**
 - Need robustness, precision and scale!!
 - « Correct-enough » and « Complete-enough » are enough (room for better definition!)
 - DSE much easier to adapt than AI
 - **New challenges and variations, so much to do**



FUTURE DIRECTION



Commissariat à l'énergie atomique et aux énergies alternatives
Institut List | CEA SACLAY NANO-INNOV | BAT. 861 – PC142
91191 Gif-sur-Yvette Cedex - FRANCE
www-list.cea.fr

Établissement public à caractère industriel et commercial | RCS Paris B 775 685 019