

Semi automatic binary deprotection

Alexandre GAZET
Yoann GUILLOT



Plan

- 1 Metasm
- 2 Structural manipulation
- 3 Challenge T2 2007
- 4 Optimization

Plan

1 Metasm

- Other disassemblers
- Binding
- Backtracing

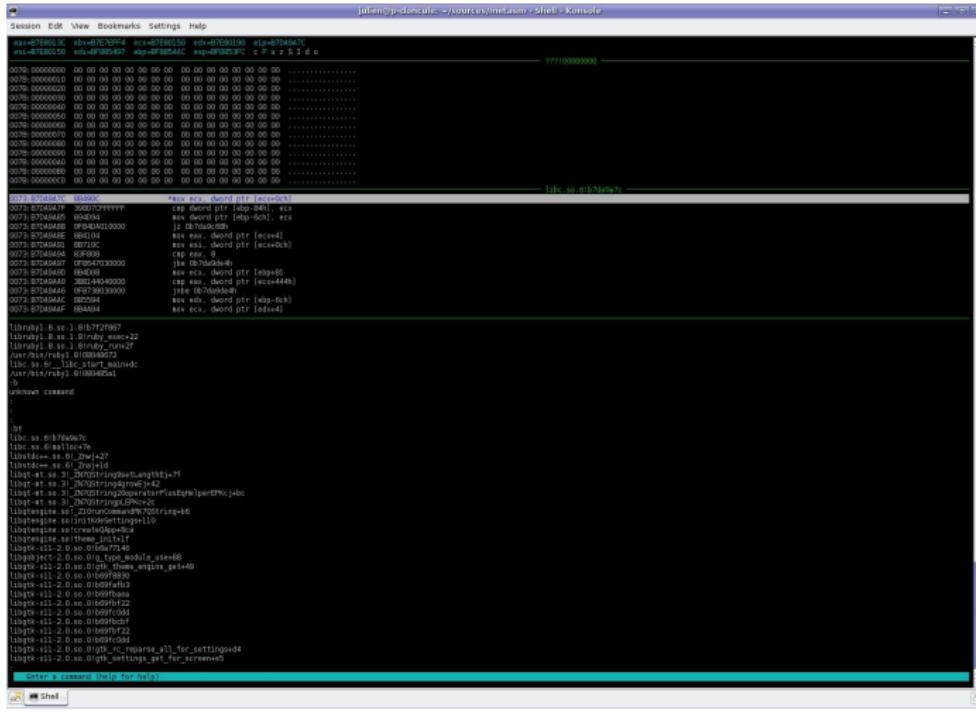
2 Structural manipulation

3 Challenge T2 2007

4 Optimization



Metasm



Metasm

The screenshot shows the Metasm disassembler interface with a control flow graph (CFG) overlaid on the assembly code. The assembly code is as follows:

```
loc_408627h:
    push 0
    push offset 412ec8
    call sub_408408h      ; x:sub_408408h
    add esp, 8
    test byte ptr [esi+10h], 1
    jz loc_4086a3h        ; x:loc_4086a3h

mov ebx, dword ptr [lat_SetFocus] ; r4:lat_SetFocus
jmp loc_408645h                 ; x:loc_408645h

// Xrefs: 408642h
loc_408645h:
    mov al, byte ptr [ebx]   ; r1:unknown
    cmp al, 20h
    js loc_408644h          ; x:loc_408644h

    cmp al, 9
    jz loc_408644h          ; x:loc_408644h

    cmp byte ptr [ebx], 22h ; r1:unknown
    jns loc_408659h         ; x:loc_408659h

// Xrefs: 408659h
loc_408659h:
    mov dl, 22h
    inc ebx
    jmp loc_40865eh         ; x:loc_40865eb

// Xrefs: 40865eh 40865bh
loc_40865eh:
    mov al, byte ptr [ebx]   ; r1:unknown
    test al, al
    jz loc_40866ch          ; x:loc_40866ch

    cmp dl, al
```

The CFG nodes are highlighted in different colors: black, grey, yellow, and red. The black nodes represent the initial entry point and the first few instructions. The grey node represents the first conditional jump. The yellow node represents the second conditional jump. The red node represents the third conditional jump. The edges between nodes are color-coded to match the nodes they connect to, showing the flow of control through the program's logic.

Disassemblers

The reference: **IDA Pro**

- Very good on *unobfuscated* code: compiled binaries (Microsoft)
- Not useful on obfuscated binaries
 - No code interpretation
 - Heavy hypothesis on code behavior

Hypothesis

- Both branches of conditionnal jumps are taken
- No overlapping instructions
- The call instruction always returns



Disassemblers

The reference: **IDA Pro**

- Very good on *unobfuscated* code: compiled binaries (Microsoft)
- Not useful on obfuscated binaries
 - No code interpretation
 - Heavy hypothesis on code behavior

Hypothesis

- Both branches of conditionnal jumps are taken
- No overlapping instructions
- The call instruction always returns



Hypothesys: call returns

```
.text:00403E9F ; -----  
.text:00403E9F  
.text:00403E9F loc_403E9F: ; CODE  
+ .text:00403E9F      push    ebp  
+ .text:00403EA0      push    ecx  
+ .text:00403EA1      push    ebp  
+ .text:00403EA2      call    sub_40BECD  
+ .text:00403EA7      outsb  
+ .text:00403EA8      cmp     edx, esp  
+ .text:00403EAA      push    esp  
+ .text:00403EAB      inc     esi  
+ .text:00403EAC      add     dword ptr [esp+4], 1  
+ .text:00403EB1      add     esp, 4  
+ .text:00403EB4      xor     ebx, edx  
+ .text:00403EB6      rep    jmp locret_4049F5  
.text:00403EBC : -----
```

```
.text:00403E9F loc_403E9F: ; CODE XREF: .text:loc_40CDEF  
.text:00403E9F      push    ebp  
.text:00403EA0      push    ecx  
.text:00403EA1      push    ebp  
.text:00403EA2      call    sub_40BECD  
.text:00403EA7      outsb  
.text:00403EA8      cmp     edx, esp  
.text:00403EAA      push    esp  
.text:00403EAB      inc     esi
```



Failure

```
push    ebp
push    ecx
push    ebp
call    sub_40BECD

db 6Eh ; n ; ====== S U B R O U T I N E =====

cmp     edx, esp
push    esp, sub_40BECD      proc near ; CODE XREF: .text:00403EA2
inc     esi
add     dword ptr [esp+0], 1
add     esp, 4
xor     ebx, edx
ret    0Ch
endp

rep    imo locret 4sub_40BECD
```

```
.text:0040BECD sub_40BECD      proc near ; CODE XREF: .text:00403EA2

.text:0040BECD
.text:0040BECF
.text:0040BED4
.text:0040BEDA

.text:0040BEDA sub_40BECD      endp
```

```
        cmp     eax, ebp
        add     dword ptr [esp+0], 1
        test   ebx, 1E2h
        ret    0Ch
```



Binding

Definition

This is how we call an instruction's semantics, through an array of symbolic expressions.

Instruction ADD:

```
a = di.instruction.args.symbolic

res = Expression [[a[0], :&, mask], :+, [a[1], :&, mask]]

binding[a[0]] = Expression[res, :&, mask]
binding[:eflag_z] = Expression[[res, :&, mask], :==, 0]
binding[:eflag_s] = sign(res)
binding[:eflag_c] = Expression[res, :>, mask]
binding[:eflag_o] = Expression[[sign(a[0]), :==, sign(a[1])],
    :'&&', [sign(a[0]), :'!=', sign(res)]]
```



Binding

Instruction CALL:

```
addr_ret = Expression[di.address, :+, di.bin_length].reduce }
binding = {
    :esp => Expression[:esp, :-, opsz],
    Indirection[:esp, opsz] => addr_ret
}
```

For exemple:

```
dword ptr [esp] = 0x4010CE
esp = esp-4
```

Instruction RDTSC:

```
binding = {
    :eax => Expression::Unknown,
    :edx => Expression::Unknown
}
```

Binding

Instruction CALL:

```
addr_ret = Expression[di.address, :+, di.bin_length].reduce }
binding = {
    :esp => Expression[:esp, :-, opsz],
    Indirection[:esp, opsz] => addr_ret
}
```

For exemple:

```
dword ptr [esp] = 0x4010CE
esp = esp-4
```

Instruction RDTSC:

```
binding = {
    :eax => Expression::Unknown,
    :edx => Expression::Unknown
}
```



Binding

Instruction CALL:

```
addr_ret = Expression[di.address, :+, di.bin_length].reduce
binding = {
    :esp => Expression[:esp, :-, opsz],
    Indirection[:esp, opsz] => addr_ret
}
```

For exemple:

```
dword ptr [esp] = 0x4010CE
esp = esp-4
```

Instruction RDTSC:

```
binding = {
    :eax => Expression::Unknown,
    :edx => Expression::Unknown
}
```



Backtracing, the theory

Definition

Symbolic emulation by walking the instruction flow backwards



Backtracing, the facts

Flot d'exécution:

```
call loc_40becdh          ; @403ea2h  e826800000
cmp eax, ebp              ; @40becdh  39e8
add dword ptr [esp+0], 1   ; @40becfh  8344240001
test ebx, 1e2h             ; @40bed4h  f7c3e2010000
ret 0ch                    ; @40bedah c20c00
```

Backtracing x dword ptr [esp] for 40bedah ret 0ch

- ❶ backtrace 40becfh add dword ptr [esp+0], 1
dword ptr [esp] => dword ptr [esp]+1
- ❷ backtrace up 40becdh->403ea2h dword ptr [esp]+1
- ❸ backtrace 403ea2h call loc_40becdh
dword ptr [esp]+1 ⇒ 403ea8h
- ❹ backtrace result: 403ea8h



Metasm

Assembler listing produced:

```
loc_403e9fh :  
    push    ebp          ; @403e9fh  55  
    push    ecx          ; @403ea0h  51  
    push    ebp          ; @403ea1h  55  
    call    loc_40becdh ; @403ea2h  e826800000  noreturn  
db 6eh                      ; @403ea7h  
  
// Xrefs: 40bedah  
loc_403ea8h :  
    cmp     edx, esp    ; @403ea8h  39e2  
    push    esp          ; @403eaah  54  
[ ... ]  
-----  
// Xrefs: 403ea2h  
loc_40becdh :  
    cmp     eax, ebp    ; @40becdh  39e8  
    add     dword ptr [esp+0], 1 ; @40becfh  8344240001  
    test    ebx, 1e2h    ; @40bed4h  f7c3e2010000  
    ret    0ch           ; @40bedah  c20c00  x: loc_403ea8h
```



Plan

1 Metasm

2 Structural manipulation

- Introduction
- Control graph complexification
- Neutral element insertion
- Unprotection

3 Challenge T2 2007

4 Optimization

Securitech 2006 - Challenge 10

Poeut.exe

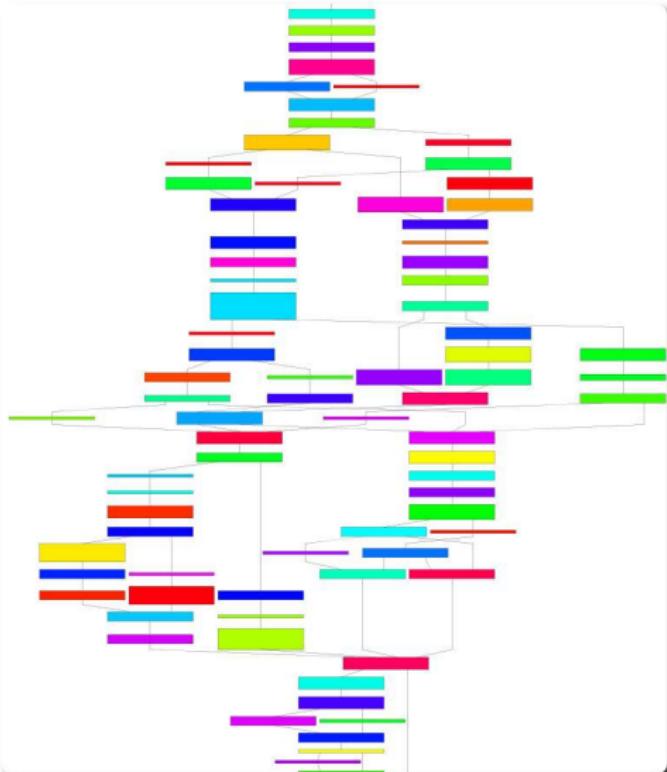
- Massively obfuscated binary
- **IDA** overwhelmed
- **Metasm** disassembles correctly, but:
 - Binary blocks are randomly moved in the binary
 - ⇒ need to write a graphic *front-end*

yEd - Graph Editor

- Visualise graphs
- Needs a *graphml* file as input
- We'll translate **Metasm** internal *InstructionBlock* representation to this format



Raw graph



Obscur predicates

Obscur predicates

- The predicate function always return true
- Using a way not easily statically analysed
- The conditionnal jump is in fact not conditionnal

```
if( x^4 * (x-5)^2 >= 0){
    goto real_code;
} else{
    goto nowhere;
}
fstp qword ptr [esp+8]
fstp qword ptr [esp]
call thunk_pow
fstp qword ptr [ebp-0x20]
mov eax, dword ptr [ebp-0ch]
sub eax, 5
push eax
fld dword ptr [esp]
lea esp, dword ptr [esp+4]
fld qword ptr [xref_8048590h]
fstp qword ptr [esp+8]
fstp qword ptr [esp]
call thunk_pow
fld qword ptr [ebp-20h]
fmulp ST(1)
```



Obscur predicates

Full randomisation

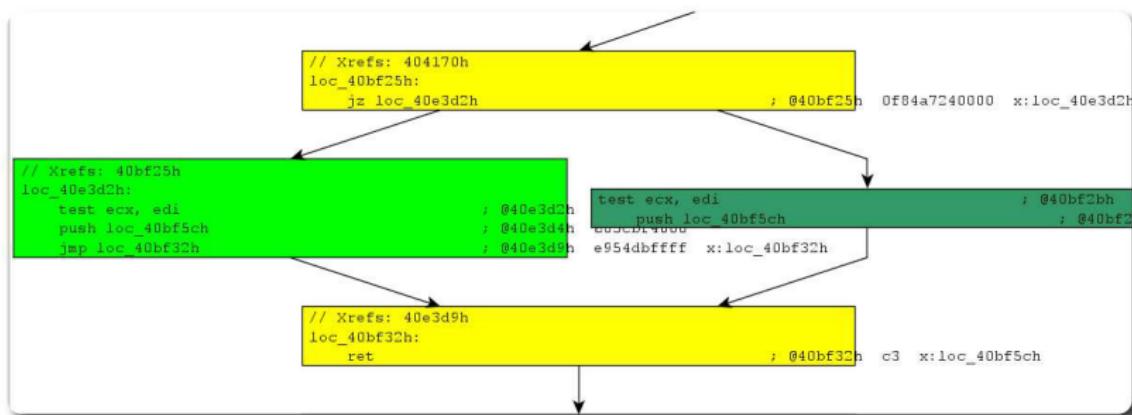
- The predicate function returns randomly true or false
- Both branches after the conditionnal jump are equivalent

```
if( rand() %2 ) {  
    real_code_A ;  
} else {  
    real_code_B ;  
}
```



Flow duplication

Massive use of random predicates:



diamond shaped graph.



Neutral elements

Definition

Instruction (or group of) having empty semantics : no effect on the execution context.



Neutral elements: apparent randomization

Implementation:

```
test esp, ebx ; @402039h 85e3
cmp ebx, edx ; @40203bh 39d3
mov byte ptr [edx], al ; @40203dh 8802
add dword ptr [esp+0], 6 ; @40203fh 8344240006
inc edx ; @402044h 42
jmp loc_40b25dh ; @402045h e913920000
```

Solving the problem

- Use the instructions' binding
- Incoherent use of the processor flags seen in the *data flow*:
 - Two successive writings,
 - or written but never read



Neutral elements: fake subfunctions

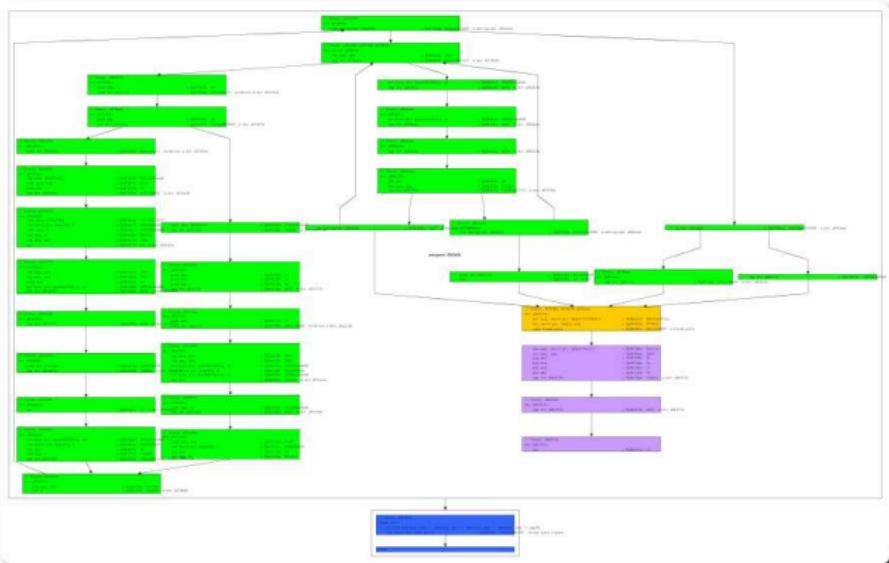
Implémentation :

```
push eax           ; @408dadh 50
push ecx           ; @408daeh 51
push ebp           ; @408dafh 55
[ ... ]
call loc_4037f2h ; @40932fh
[ ... ]
push esp           ; @4037f4h 54
push ecx           ; @4037f5h 51
[ ... ]
add dword ptr [esp+8], 9 ; @4037f9h 8344240809
add esp, 8         ; @403800h 83c408
[ ... ]
ret 0ch            ; @403808h c20c00 x:loc_40933dh
```

Solution

- Execution flow reconstruction
- Stack emulation
- We have a pattern: return address modification

Epilogue (raw)



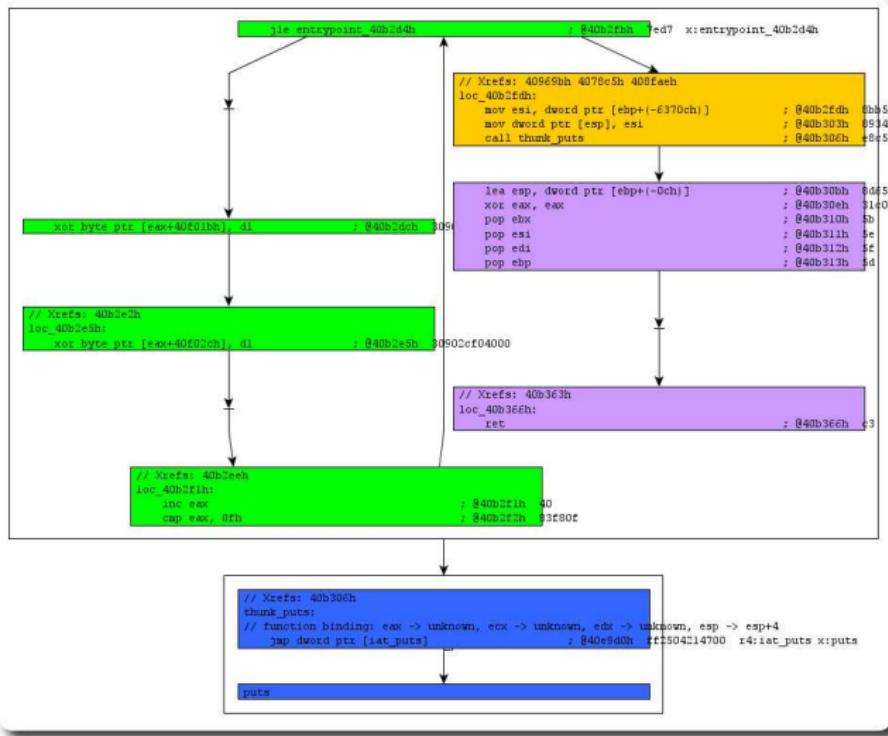
Execution flow analysis

Solution

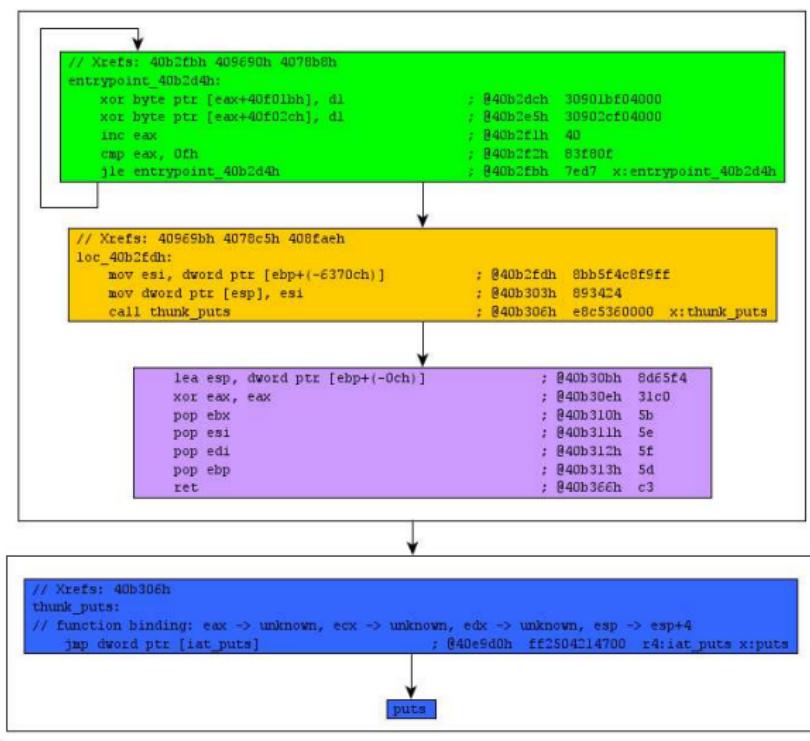
- ➊ Walk the internal tree of InstructionBlock
- ➋ Inline functions if needed
- ➌ Scan for a diamond pattern
- ➍ Construction, cleaning and comparison of flows
- ➎ Factorisation



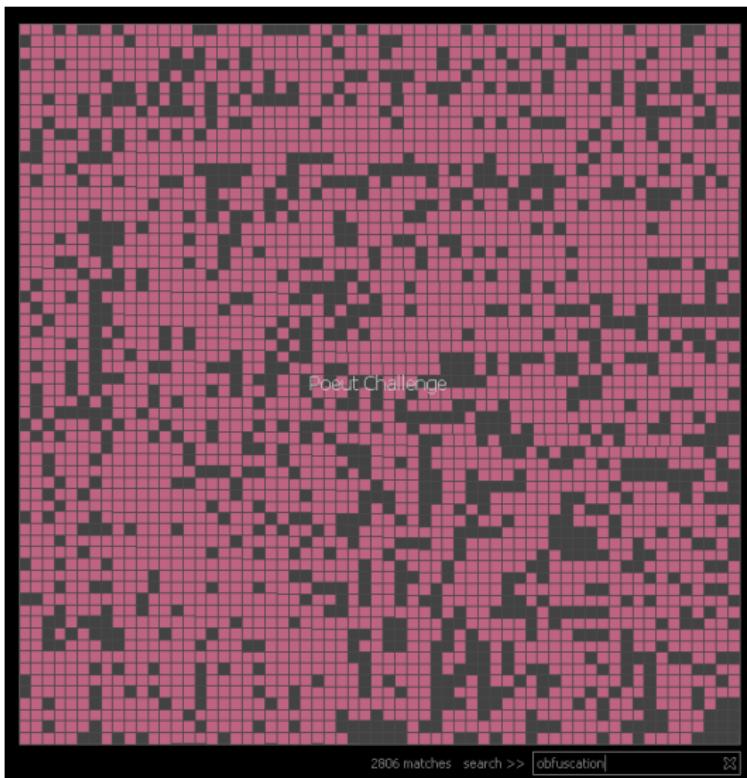
Epilogue (factorized)



Epilogue (reordered)



Epilogue (final)



2806 matches search >> obfuscation



Last touch

- Cleanup of the whole program graph
- Output of a clean asm source
- Used to reassemble an unprotected binary



Plan

- 1 Metasm
- 2 Structural manipulation
- 3 Challenge T2 2007
 - Introduction
 - Obfuscation
 - Virtual machine
 - Resolution
- 4 Optimization



T2 2007

t207.exe

- <http://www.t2.fi/challenge/>
- Goal: find the password to unlock the program
- Very simple binary [demo/1]
- Loads an obfuscated driver [demo/2]



Deobfuscation

Obfuscation types

- Junk code
- Obfuscated arithmetics
- Ring3 detection
- Code duplication



Junk code

junk

```
ror edi, 0dh
xchg ebx, edi
ror ebx, 13h
xchg ebx, edi
```

Obfuscated arithmetics

bit rotation

```
push eax
push ecx
rol dword ptr [esp+4], cl
pop ecx
pop eax
```



Ring3 detection

test ring0

```
pushfd
push eax
xor eax, eax
mov ax, cs
cmp eax, 9
jle loc_131d5h
rdtsc
imul eax, ecx
jmp eax ; x:unknown
```

loc_131d5h:

```
pop eax
popfd
```

Code duplication

duplication

```
push esi
push ebx
pushfd
rdtsc
imul ecx, ebx
cmp cl, 7fh
jnb loc_21aba
popfd
pop ebx
pop esi

loc_21abah:
popfd
pop ebx
pop esi
```



Virtual machine

Handler structure

- Simple operations
- Similar blueprint
- Operations controlled by [ebp]

[demo/3]



Virtual addition

An addition handler

```
loc_15336:  
    mov ecx, dword ptr [ebp+0ch]  
    xor ecx, 842b1208h  
    mov ecx, dword ptr [ebx+ecx]  
    mov eax, dword ptr [ebp+8]  
    xor eax, 842b1208h  
    add dword ptr [ebx+eax], ecx
```

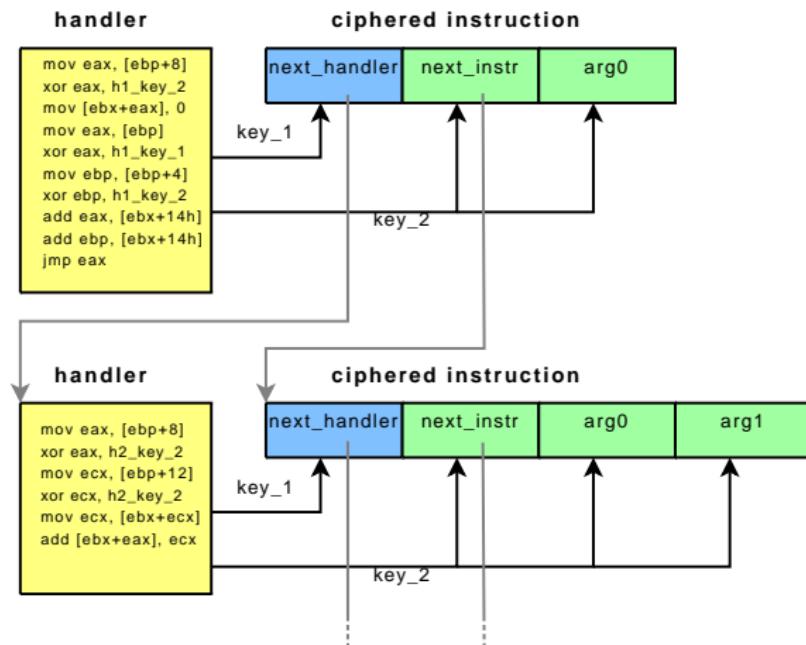
Inter-handler transition

transition to the next handler

```
mov ecx, dword ptr [ebp+0]
xor ecx, 149f0c63h
mov ebp, dword ptr [ebp+4]
xor ebp, 842b1208h
add ebp, dword ptr [ebx+14h]
add ecx, dword ptr [ebx+14h]
jmp ecx      ; x:unknown
```



Virtual machine architecture



Handler enumeration

Direct enumeration impossible

- No “handler table”
- Need to follow the virtual code flow step by step
- Binary size implies numerous handlers
- Auto-analyse every handler behavior using backtracking

[demo/4]



Analyse result

Handler binding

```
handler_13491h:  
// "reg0" <- Expression["reg0", :+, "reg1"]  
// handler type: add reg, reg  
mov eax, dword ptr [ebp+0ch]  
xor eax, 8d3f5d8bh  
mov eax, dword ptr [ebx+eax]  
mov ecx, dword ptr [ebp+8]  
xor ecx, 8d3f5d8bh  
add dword ptr [ebx+ecx], eax
```



Raw virtual code

The first virtual instructions

```
entrypoint_219feh_21ea6h:  
    nop  
    mov r68, 28h  
    add r68, host_esp  
    mov r64, dword ptr [r68]  
    mov dword ptr [esp], r64  
    mov r64, 4  
    add esp, r64  
    mov r68, 2ch  
    add r68, host_esp  
    mov r64, dword ptr [r68]  
    mov dword ptr [esp], r64  
    mov r64, 4  
    add esp, r64  
    trap
```



Virtual macro-instructions

Higher abstraction level

- Through pattern recognition
- Reconstruction of a higher level assembler
- Apparition of functions (call, ret)



Virtual macro-code

The first instructions (again)

```
entrypoint_219feh_21ea6h:  
    mov dword ptr [esp], dword ptr [host_esp+28h]  
    add esp, 4  
    mov dword ptr [esp], dword ptr [host_esp+2ch]  
    add esp, 4  
    mov ebp, esp  
    add esp, 234h  
    mov r64, dword ptr [ebp+200h]  
    xor r64, 1  
    jrz loc_2d630h_2d8ffh, r64  
    syscall_alloc_ptr r64, 0ch  
    mov dword ptr [ebp+200h], r64  
loc_2d630h_2d8ffh:
```



Decompilation

- The instruction semantic is very simple
- The code patterns reminds C
- Check it out

Validation

- Enabled us to solve the challenge by a pure static approach



Plan

1 Metasm

2 Structural manipulation

3 Challenge T2 2007

4 Optimization

- Introduction
- POC
- Results
- Prospects

Facts

What we have done so far

- Stress was put on binary manipulation : **code rewriting**
- Add effective methods in Metasm to clean obfuscation on the fly
- At various level: filtering processor, graph manipulation ...

Drawbacks

- Lack of a higher level of abstraction
- Analysis is still time-consuming and painful
- Mostly target-specific



Concepts

Objectives

- Generic rewriting rules
- Obfuscated code as input
- Clean (obfuscation free) code as output

Means

- What we want to do looks like optimization !
- Extensive litterature on the subject
- Can be easily applied at assembly level using Metasm: we have methods to work on basic blocks & instructions



Proof of concept

Implemented optimizations

- **Declaration cleaning:** remove useless assignments
- **Constant propagation**
- **Constant, operation folding:** apply basic rules of arithmetic
- **Peephole:** replace known patterns with a reduced form

Each of those optimization amounts to one or many rewriting rules,
possibly associated with a condition.

We apply them locally, on each basic bloc, spaghetti code has been
previously merged.s



Example: constant propagation

Before

```
100bb5cfh mov al, 12h
100bed67h mov cl, 46h
100bed69h xor cl, al
```

After

```
100bb5cfh mov al, 12h
100bed67h mov cl, 46h
100bed69h xor cl, 12h
```

The constant (0x12) has been propagated through *al*



Example 2: constant folding

Before

```
100bb5cfh mov al, 12h
100bed67h mov cl, 46h
100bed69h xor cl, 12h
```

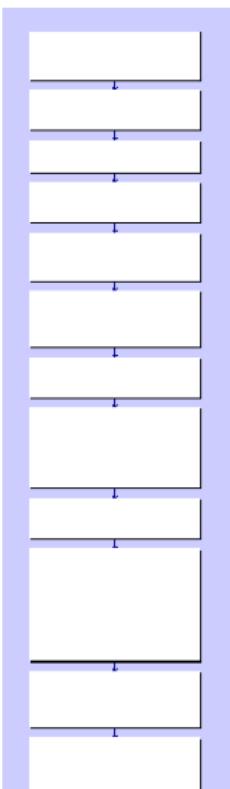
After

```
100bb5cfh mov al, 12h
100bed67h mov cl, 54h
```

cl is now assigned with $0x46 \text{ xor } 0x12h$



Does it work?



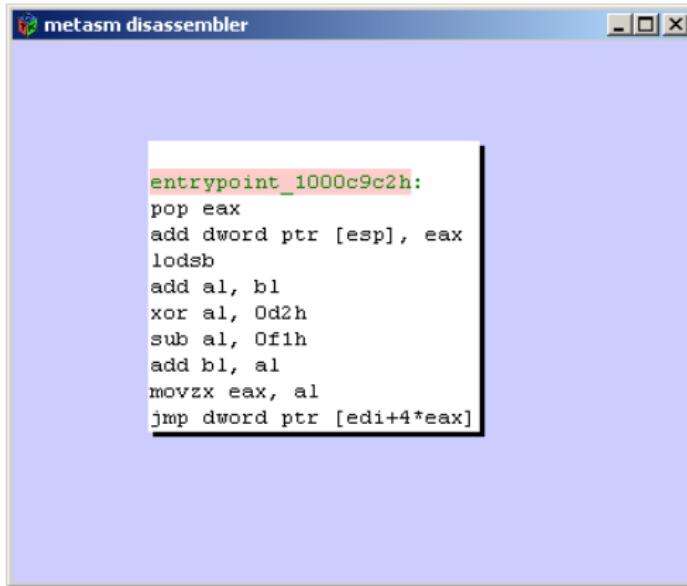
Does it work ?



The screenshot shows the Metasm disassembler interface with several assembly code snippets. Each snippet is preceded by a green double-headed arrow indicating a structural manipulation operation. The snippets are:

- `entrypoint_1000c9c2h:
push dword ptr [esp]
pop eax
jmp loc_1000e27eh ; x:loc_1000e27eh`
- `// Xrefs: 1000c9c6h
loc_1000e27eh:
push ebx
mov ebx, esp
add ebx, 4
add ebx, 4
xchg ebx, dword ptr [esp]
pop esp
add dword ptr [esp], 4af11ef6h
add dword ptr [esp], eax
push ebx
push 4af11ef6h
pop ebx
sub dword ptr [esp+4], ebx
jmp loc_1000c76dh ; x:loc_1000c76dh`
- `// Xrefs: 1000e2a3h
loc_1000c76dh:
pop ebx
jmp loc_1000ec16h ; x:loc_1000ec16h`
- `// Xrefs: 1000c76eh
loc_1000ec16h:
jmp loc_1000bb4eh ; x:loc_1000bb4eh`
- `// Xrefs: 1000ec16h
loc_1000bb4eh:
lodsb
jmp loc_1000d4ech ; x:loc_1000d4ech`
- `// Xrefs: 1000bb4fh
loc_1000d4ech:
add al, -7fh
add al, -70h`

Does it work ? Actually yes



The screenshot shows the Metasm disassembler interface. The window title is "metasm disassembler". The assembly code is displayed in a text area:

```
entrypoint_10000c9c2h:  
pop eax  
add dword ptr [esp], eax  
lodsb  
add al, bl  
xor al, 0d2h  
sub al, 0f1h  
add bl, al  
movzx eax, al  
jmp dword ptr [edi+4*eax]
```



Conclusion and prospect

Conclusion

- Using optimization to defeat obfuscation is very promising
- The kind of obfuscation used in the protection is too weak: mainly based on local constants expansion and affine functions.
- We don't have a real intermediate representation

Prospect

- We are working on decompilation problems
- We will study the opportunity to use jointly **Metasm** and **LLVM**



The end

Thank you for your attention !

Any questions ?

