**Reverse engineering through execution trace analysis**

Gabriel Duque & Tanguy Dubroca

LSE

- Reversing virtual machines is a hassle
- Their state is optimizable

- Capture traces with different inputs
- Analyze traces to build a CFG
- Lift our CFG to LLVM bytecode
- Run several optimization passes
- Rebuild an executable or analyze the LLVM IR

**Figure 1:** trace

```
0x1337a707: mov dword ptr [rbx - 4], eax
r14 : 0x0000000000000000 r15 : 0x000000009b7fdbf0 rcx : 0x000000001337a707 rsi : 0x00007f39fbaa4fb8
r10 : 0x0000000000000008 rbx : 0x00007f39fba9400c rsp : 0x00007ffeae7a3708 r11 : 0x0000000000000246
 r8 : 0x0000000000000000 rdx : 0x0000000000000020 rip : 0x000000001337a707  r9 : 0x00007f39fba72be0
r12 : 0x000000000000ff8c rbp : 0x00007f39fba95000 rdi : 0x0000000000000000 rax : 0x0000000000000020
-----------------------------------------
0x1337a70a: sub rbx, 4
r14 : 0x0000000000000000 r15 : 0x000000009b7fdbf0 rcx : 0x000000001337a707 rsi : 0x00007f39fbaa4fb8
r10 : 0x0000000000000008 rbx : 0x00007f39fba9400c rsp : 0x00007ffeae7a3708 r11 : 0x0000000000000246
 r8 : 0x0000000000000000 rdx : 0x0000000000000020 rip : 0x000000001337a70a  r9 : 0x00007f39fba72be0
r12 : 0x000000000000ff8c rbp : 0x00007f39fba95000 rdi : 0x0000000000000000 rax : 0x0000000000000020
-----------------------------------------
0x1337a70e: jmp 0x133707de
r14 : 0x0000000000000000 r15 : 0x000000009b7fdbf0 rcx : 0x000000001337a707 rsi : 0x00007f39fbaa4fb8
r10 : 0x0000000000000008 rbx : 0x00007f39fba94008 rsp : 0x00007ffeae7a3708 r11 : 0x0000000000000246
 r8 : 0x0000000000000000 rdx : 0x0000000000000020 rip : 0x000000001337a70e  r9 : 0x00007f39fba72be0
r12 : 0x000000000000ff8c rbp : 0x00007f39fba95000 rdi : 0x0000000000000000 rax : 0x0000000000000020
...
```

- Traces are raw data
- They are not constrained to a specific use case
- Example projects:
  - Griffin: https://www.microsoft.com/en-us/research/wp-content/uploads/2017/01/griffin-asplos17.pdf
  - HeNet: https://arxiv.org/pdf/1801.02318

- Syscall giving control over a process

- Can singlestep through the code

- Gives access to registers

- **But** painfully slow

- Allows to debug any kind of code

- Quite easy to implement

- Extensible

- **But** may not execute correctly the target executables

- Frameworks such as DynamoRIO, Pintools, Valgrind

- Very fast compared to previous methods

- **But** may add significant overhead during execution

- In our case Intel PT

- Can trace anything (userland, kernel, hypervisor)

- Fastest method

- Low overhead (around 5% for intel pt)

- **But** difficult setup, trace loss and decoding overhead

- ptrace
- intel PT

| size | Description |
| --- | --- |
| 8 | Magic (0xe9cae282c414b97d) |
| 8 | Edge count |
| edge count * sizeof(trace_entry) | Edge entries |
| 8 | Program mapping count |
| mapping count * sizeof(mapping_entry) | Memory mapping entries |

# Trace format



**Figure 2:** ghidra plugin

Necessary information for native code lifting:

- Functions

- Basic blocks

    - Succesors

- Instructions

    - Cross-references
        - Code
        - Data

Two passes:

- First pass: build a linear CFG
  - Iterate through instructions
  - Add them to the current basic block
  - Start a new basic block when we find a control flow instruction
- Second pass: deduce the real CFG
  - Detect loops
  - Merge parents and children vectors

id: 1
begin: 0x100001149
end: 0x10000114d
instr count: 2

0x100001149: cmpl $0, -4(%rbp)
0x10000114d: jne 0xffffffffffffffec

id: 2
begin: 0x100001139
end: 0x100001144
instr count: 5

0x100001139: movl -4(%rbp), %eax
0x10000113c: leal -1(%rax), %edx
0x10000113f: movl %edx, -4(%rbp)
0x100001142: movl %eax, %edi
0x100001144: callq 0xffffffffffffffd5

id: 3
begin: 0x100001119
end: 0x100001127
instr count: 7

0x100001119: pushq %rbp
0x10000111a: movq %rsp, %rbp
0x10000111d: movl %edi, -4(%rbp)
0x100001120: movl -4(%rbp), %eax
0x100001123: addl $2, %eax
0x100001126: popq %rbp
0x100001127: retq

id: 4
begin: 0x100001149
end: 0x10000114d
instr count: 2

0x100001149: cmpl $0, -4(%rbp)
0x10000114d: jne 0xffffffffffffffec

id: 5
begin: 0x100001139
end: 0x100001144
instr count: 5

0x100001139: movl -4(%rbp), %eax
0x10000113c: leal -1(%rax), %edx
0x10000113f: movl %edx, -4(%rbp)
0x100001142: movl %eax, %edi
0x100001144: callq 0xffffffffffffffd5

**Figure 3:** linear

**Figure 4:** color

**Figure 5:** loop

- Disassemble with capstone

- Analyze the flow of the binary

- Generate the CFG

**Figure 6:** nucleus 0

- Ignore *call* edges
- Basic blocks are connected through intraprocedural edges
- Detect basic block clusters

**Figure 7:** nucleus 1

- Reintroduce *call* edges

- Start at entrypoints

- Follow recursively until complete functions are formed

**Figure 8:** nucleus 2

We have:

- Functions

We have:

- Functions
- Blocks

We have:

- Functions
- Blocks
- Instructions

**Cross-references**

- Internal or external ?

- Code or data ?

- Lifting works.
- Optimizing works.
- Rebuilding works.

**Figure 9:** meme1

```c
  1
  2  void sub_40076d_main(astruct *param_1,long param_2)
  3
  4  {
  5    int *piVar1;
  6    int *piVar2;
  7    int iVar3;
  8    long *plVar4;
  9    undefined8 uVar5;
 10    undefined *puVar6;
 11    int *piVar7;
 12    byte bVar8;
 13    uint uVar9;
 14    byte bVar10;
 15    uint uVar11;
 16    int iVar12;
 17    long lVar13;
 18    uint uVar14;
 19    uint uVar15;
 20    long lVar16;
 21    ulong uVar17;
 22    uint uVar18;
 23    long lVar19;
 24    bool bVar20;
 25    bool bVar21;
 26
 27    lVar16 = param_1->field_0x908;
 28    *(long *)(lVar16 + -8) = param_1->field_0x918;
 29    param_1->field_0x918 = lVar16 + -8;
 30    param_1->field_0x908 = lVar16 + -0x68;
 31    *(undefined4 *)(lVar16 + -0x5c) = *(undefined4 *)&param_1->field_0x8f8;
 32    *(ulong *)(lVar16 + -0x68) = param_1->field_0x8e0;
 33    *(undefined8 *)(lVar16 + -0x10) = *(undefined8 *)(param_1->field_0x878 + 0x28);
 34    *(undefined4 *)(lVar16 + -0x44) = 0;
 35    *(undefined8 *)(lVar16 + -0x4c) = 0x100000001;
 36    iVar12 = *(int *)(lVar16 + -0x4c);
 37    param_1->field_0x8c8 = (long)iVar12;
 38    uVar17 = (long)*(int *)(lVar16 + -0x48) * 0xb + (long)iVar12;
 39    puVar6 = maze + uVar17;
 40    param_1->field_0x811 = -(-(puVar6 < maze) | -(0xffffffffff9fbf7f < uVar17)) & 1;
 41    uVar14 = (int)((ulong)puVar6 & 0xff) - ((uint)((ulong)puVar6 & 0xff) >> 1) & 0x55);
 42    uVar14 = (uVar14 >> 2 & 0x33333333) + (uVar14 & 0x33333333);
 43    param_1->field_0x813 = -(byte)(((uVar14 >> 4) + uVar14 & 0x0f0f0f0f) * 0x1010101 >> 0x18) & 1;
 44    param_1->field_0x815 = ((byte)*uVar17 ^ 0x80 ^ (byte)puVar6) >> 4 & 1;
 45    param_1->field_0x817 = puVar6 == (undefined *)0x0;
 46    param_1->field_0x819 = (byte)((ulong)puVar6 >> 0x3f);
 47    param_1->field_0x81d = ((ulong)puVar6 >> 0x3f) + ((uVar17 ^ (ulong)puVar6) >> 0x3f) == 2;
 48    *puVar6 = 0x58;
 49    param_1->field_0x8b8 = 7;
 50    param_1->field_0x8e0 = 0xb;
 51    param_1->field_0x8f8 = 0x403518;
 52    param_1->field_0x8a8 = 0;
 53    lVar16 = param_1->field_0x908;
 54    *(long *)(lVar16 + -8) = param_2 + 0x74;
 55    param_1->field_0x908 = lVar16 + -8;
 56    param_1->field_0x9a8 = param_2 + -0x1cd;
 57    uVar5 = ext_6010e8_printf(param_1);
 58    lVar16 = param_1->field_0x9a8;
 59    param_1->field_0x8d8 = (ulong)*(uint *)(param_1->field_0x918 + -0x40);
 60    param_1->field_0x8e8 = (ulong)*(uint *)(param_1->field_0x918 + -0x44);
 61    param_1->field_0x8f8 = 0x403530;
 62    param_1->field_0x8a8 = 0;
 63    lVar19 = param_1->field_0x908;
 64    *(long *)(lVar19 + -8) = lVar16 + 0x17;
 65    param_1->field_0x908 = lVar19 + -8;
 66    param_1->field_0x9a8 = lVar16 + -0x24l;
 67    uVar5 = ext_6010e8_printf(param_1,uVar5);
 68    lVar16 = param_1->field_0x9a8;
 69    param_1->field_0x8e8 = (ulong)*(uint *)(param_1->field_0x918 + -0x3c);
 70    param_1->field_0x8f8 = 0x403548;
 71    param_1->field_0x8a8 = 0;
 72    lVar19 = param_1->field_0x908;
 73    *(long *)(lVar19 + -8) = lVar16 + 0x14;
 74    param_1->field_0x908 = lVar19 + -8;
 75    param_1->field_0x9a8 = lVar16 + -600;
 76    uVar5 = ext_6010e8_printf(param_1,uVar5);
 77    lVar16 = param_1->field_0x9a8;
 78    param_1->field_0x8f8 = 0x403560;
 79    lVar19 = param_1->field_0x908;
 80    *(long *)(lVar19 + -8) = lVar16 + 10;
 81    param_1->field_0x908 = lVar19 + -8;
 82    param_1->field_0x9a8 = lVar16 + -0x28c;
 83    uVar5 = ext_6010d8_puts(param_1,uVar5);
 84    lVar16 = param_1->field_0x9a8;
 85    param_1->field_0x8f8 = 0x4035a2;
 86    lVar19 = param_1->field_0x908;
 87    *(long *)(lVar19 + -8) = lVar16 + 10;
 88    param_1->field_0x908 = lVar19 + -8;
 89    param_1->field_0x9a8 = lVar16 + -0x296;
 90    uVar5 = ext_6010d8_puts(param_1,uVar5);
 91    lVar16 = param_1->field_0x908;
 92    lVar19 = param_1->field_0x908;
 93    *(long *)(lVar19 + -8) = lVar19 + 5;
 94    param_1->field_0x908 = lVar16 + -8;
 95    sub_4006f6_draw(param_1,lVar19 + -0x12a,uVar5);
 96    lVar16 = param_1->field_0x9a8;
 97    uVar17 = param_1->field_0x918 - 0x30;
 98    param_1->field_0x8a8 = uVar17;
 99    param_1->field_0x8d8 = 0x1c;
100    param_1->field_0x8e8 = uVar17;
101    param_1->field_0x8f8 = 0;
102    lVar19 = param_1->field_0x908;
103    *(long *)(lVar19 + -8) = lVar16 + 0x16;
104    param_1->field_0x908 = lVar19 + -8;
105    param_1->field_0x9a8 = lVar16 + -0x275;
106    uVar5 = ext_6010f0_read(param_1);
107    lVar16 = param_1->field_0x8a8 + 0x1ee;
108    do {
109      lVar19 = param_1->field_0x918;
110      uVar15 = *(uint *)(lVar19 + -0x3c);
111      uVar14 = uVar15 - 0x1b;
112      *(bool *)&param_1->field_0x811 = uVar15 < 0x1b;
```
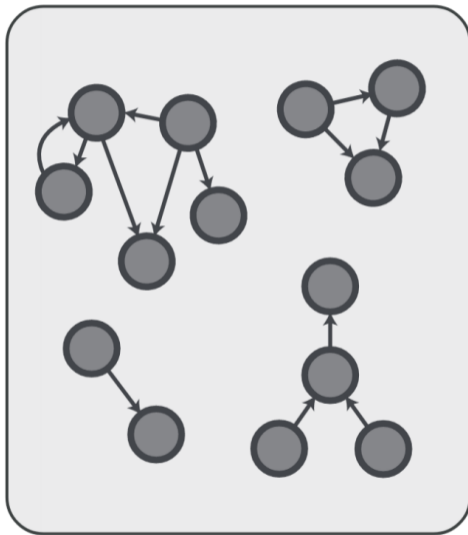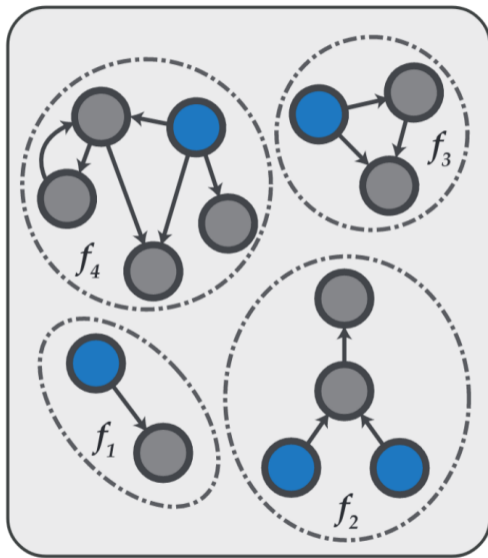
```
111    uVar14 = uVar15 - 0x1b;
112    *(bool *)&param_1->field_0x811 = uVar15 < 0x1b;
113    uVar11 = (uVar14 & 0xff) - (uVar14 >> 1 & 0x55);
114    uVar11 = (uVar11 >> 2 & 0x33333333) + (uVar11 & 0x33333333);
115    param_1->field_0x813 = ~(byte)(((uVar11) >> 4) + uVar11 & 0x0f0f0f0f) * 0x1010101 >> 0x18) & 1;
116    param_1->field_0x815 = ~(byte)(uVar14 ^ uVar15) >> 4 & 1;
117    bVar20 = uVar14 == 0;
118    param_1->field_0x817 = bVar20;
119    bVar8 = (byte)(uVar14 >> 0x1f);
120    param_1->field_0x819 = bVar8;
121    bVar21 = ((uVar14 ^ uVar15) >> 0x1f) + (uVar15 >> 0x1f) == 2;
122    param_1->field_0x81d = bVar21;
123    bVar21 = (bVar8 != 0) == bVar21;
124    lVar13 = 10;
125    if (bVar20 || !bVar21) {
126      lVar13 = -0x1e9;
127    }
128    lVar13 = lVar13 + lVar16;
129    if (bVar21 && !bVar20) {
130      param_1->field_0x8f8 = 0x403642;
131      lVar16 = param_1->field_0x908;
132      *(long *)(lVar16 + -8) = lVar13 + 10;
133      param_1->field_0x908 = lVar16 + -8;
134      param_1->field_0x9a8 = lVar13 + -0x4b3;
135      uVar5 = ext_6010d8_puts(param_1,uVar5);
136      param_1->field_0x8a8 = 1;
137      plVar4 = (long *)param_1->field_0x918;
138      uVar17 = *(ulong *)(param_1->field_0x878 + 0x28) ^ plVar4[-1];
139      param_1->field_0x8e8 = uVar17;
140      param_1->field_0x811 = 0;
141      uVar14 = (int)(uVar17 & 0xff) - ((uint)((uVar17 & 0xff) >> 1) & 0x55);
142      uVar14 = (uVar14 >> 2 & 0x33333333) + (uVar14 & 0x33333333);
143      param_1->field_0x813 = ~(byte)(((uVar14 >> 4) + uVar14 & 0x0f0f0f0f) * 0x1010101 >> 0x18) & 1;
144      bVar21 = uVar17 == 0;
145      param_1->field_0x817 = bVar21;
146      param_1->field_0x819 = (byte)(uVar17 >> 0x3f);
147      param_1->field_0x81d = false;
148      param_1->field_0x815 = 0;
149      if (bVar21) {
150        param_1->field_0x918 = *plVar4;
151        param_1->field_0x9a8 = plVar4[1];
152        *(long **)&param_1->field_0x908 = plVar4 + 2;
153        return;
154      }
155      lVar16 = param_1->field_0x9a8 + 0x14 + (ulong)bVar21 * 5;
156      lVar19 = param_1->field_0x908;
157      *(long *)(lVar19 + -8) = lVar16 + 5;
158      param_1->field_0x908 = lVar19 + -8;
159      param_1->field_0x9a8 = lVar16 + -0x4c1;
160      uVar5 = ext_601000___stack_chk_fail(param_1,lVar5);
161      goto LAB_00401ad5;
162    }
163    *(undefined4 *)(lVar19 + -0x38) = *(undefined4 *)(lVar19 + -0x44);
164    *(undefined4 *)(lVar19 + -0x34) = *(undefined4 *)(lVar19 + -0x40);
165    bVar8 = *(byte *)(lVar19 + -0x30) + (long)*(int *)(lVar19 + -0x3c));
166    uVar14 = SEXT14((char)bVar9);
167    piVar1 = (int *)(lVar19 + -0x44);
```

**Figure 10:** cancer 3

**Figure 11:** meme3

**Wait a minute... Was this supposed to be a deobfuscation or an obfuscation tool?**

# Custom LLVM IR wrapper

```c
#include <stddef.h>

struct State {                    static void sub_loop_0(struct State *s_ptr) {
    size_t rdi;                       s_ptr->rax = s_ptr->rdi;
    size_t rsi;                       for (s_ptr->r11 = 0; s_ptr->r11 != s_ptr->rsi; ++(s_ptr->r11))
    size_t rax;                           s_ptr->rax += s_ptr->rsi;
    size_t r11;                       return;
};                                }

size_t wrapper_loop(size_t a, size_t b) {
    struct State s;
    s.rdi = a;
    s.rsi = b;

    sub_loop_0(&s);
    return s.rax;
}
```

```
define i64 @wrapper_loop(i64, i64) {
    %3 = alloca %struct.State
    ; set struct.State.rdi and struct.State.rsi before entering @sub_loop_0
    ; %4 is a pointer to struct.State.rdi
    ; %5 is a pointer to struct.State.rsi
    %4 = getelementptr %struct.State, %struct.State* %3, i32 0, i32 0
    store i64 %0, i64* %4
    %5 = getelementptr %struct.State, %struct.State* %3, i32 0, i32 1
    store i64 %1, i64* %5

    ; call our static sub function
    call void @sub_loop_0(%struct.State* %3)

    ; return the right value: struct.State.rax
    %6 = getelementptr %struct.State, %struct.State* %3, i32 0, i32 2
    %7 = load i64, i64* %6
    ret i64 %7
}
```

**Figure 12:** meme2

A big problem still remains memory access generates function calls to intrinsic procedures.

Because of these generated calls to functions in *remill* functions are not optimizable.

We have to write an optimization pass to force-inline them.

- Looked like a good idea.

- Was not a good idea.

- We would have to rewrite a big part of *mcsema*/*remill*.

EPITA

Do you have any questions?