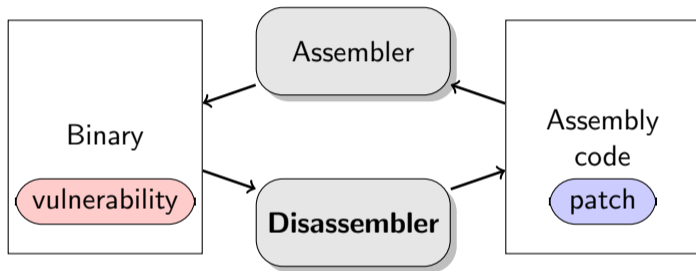


Sometimes the **source code is not available**



Sometimes the **source code is not available**

Disassembly is the first step in:

- ▶ Binary analysis
- ▶ Binary rewriting/hardening

We want to obtain *reassemblable assembly*

- ▶ Assembly code with cross references
- ▶ We can modify it without breaking it

This involves two tasks:

- ▶ Instruction Boundary Identification
- ▶ Symbolization

Instruction Boundary Identification

A binary looks like this:

```
ba 08 3b 44 00 31 c9 0f  
1f 80 00 00 00 00 48 39  
3a ... 48 83 c2 20 48 83  
f9 3a 7e ee
```

Instruction Boundary Identification

A binary looks like this:

```
ba 08 3b 44 00 | 31 c9 | 0f
1f 80 00 00 00 00 | 48 39
3a | ... | 48 83 c2 20 | 48
83 f9 3a | 7e ee |
```



```
40c7f2:  mov EDX, 443b08
40c7f7:  xor ECX, ECX
40c7f9:  nop
40c800:  cmp [RDX], RDI
:      :
40c808:  add RDX, 32
40c80c:  cmp RCX, 58
40c810:  jle 40c800
```

Instruction Boundary Identification amounts to finding which addresses correspond to which instructions.

Instruction Boundary Identification

Challenging because:

- ▶ X64 instructions have variable sizes
- ▶ Data interleaves with instructions

Symbolization: Which numbers are references vs. literals

```
40c7f2: mov EDX, 443b08
40c7f7: xor ECX,ECX
40c7f9: nop
40c800: cmp [RDX],RDI
:
:
40c808: add RDX,32
40c80c: cmp RCX,58
40c810: jle 40c800
:
:
```



```
mov EDX, Label_1
xor ECX,ECX
nop
Label_2: cmp [RDX],RDI
:
:
add RDX,32
cmp RCX,58
jle Label_2
:
```

```
443b40: 04 04 00 00 00 00 00 00
443b48: d0 50 41 00 00 00 00 00
443b50: 85 48 44 00 00 00 00 00
```

```
04 04 00 00 00 00 00 00
.quad Label_3
.quad Label_4
```

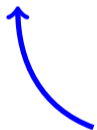

Symbolization: Which numbers are references vs. literals

```
40c7f2:  mov EDX, 443b08
40c7f7:  xor ECX,ECX
40c7f9:  nop
40c800:  cmp [RDX],RDI
```

Symbolic references allow us to move and modify the code without breaking it

```
443b40:  04 04 00 00 00 00 00 00
443b48:  d0 50 41 00 00 00 00 00
443b50:  85 48 44 00 00 00 00 00
```

```
mov EDX, Label_1
xor ECX,ECX
nop
Label_2:  cmp [RDX],RDI
:
add RDX,32
cmp RCX,58
jle Label_2
:
04 04 00 00 00 00 00 00
.quad Label_3
.quad Label_4
```



Symbolization: Which numbers are references vs. literals

40c7f2: mov EDX, 443b08

40c7f7: xor ECX,ECX

40c7f9: nop

40c800: cmp [RDX],RDI

mov EDX, **Label_1**

xor ECX,ECX

nop

Label_2: cmp [RDX],RDI

Symbolic references are
and modify the code with

Relocation information is
NOT enough (even for PIE)

443b40: 04 04 00 00 00 00 00 00

443b48: d0 50 41 00 00 00 00 00

443b50: 85 48 44 00 00 00 00 00

04 04 00 00 00 00 00 00

.quad **Label_3**

.quad **Label_4**

Our approach: Use Datalog

Both **instruction boundary identification** and **symbolization** are hard problems

We want to:

- ▶ Combine different heuristics easily
- ▶ Use simple static analyses to inform our decisions
- ▶ Run it quickly

Binaries are encoded as **facts** (our initial knowledge of the binary)

- ▶ We decode **every possible offset** in code sections (obtaining a **superset** of all possible instructions in the code)

Instruction:

```
4000A0: mov RAX, 420020
```

Becomes:

```
instruction(4000A0,5,' ','MOV',1,2,0,0)
op_regdirect(1,'RAX')
op_immediate(2,420020)
```

If the decoding fails at address **A**, we generate `invalid(A)`

Analyses and heuristics are expressed as Datalog **rules**:

Backward traversal that propagates invalid instructions

```
invalid(From):-  
  (  
    must_fallthrough(From,To);  
    direct_jump(From,To);  
    direct_call(From,To)  
  ),(  
    invalid(To);  
    !instruction(To,_,_,_,_,_,_,_)  
  ).
```

If there is an instruction at address **From** that must fallthrough, or jumps, or calls an address **To** that contains an invalid instruction or no instruction at all, then the instruction at **From** is also invalid.

Instruction Boundary Identification

1. Backward traversal: propagate invalid instructions
2. Forward traversal: build **superset** of all possible basic blocks
 - ▷ Hybrid between linear sweep and recursive traversal
 - ▷ Use potential references in data sections
3. Assign points to candidate blocks using heuristics
 - ▷ Entry point: +20
 - ▷ Address appears in data section: +1
 - ▷ Direct jump: +6
 - ▷ ...
4. Aggregate points to resolve overlaps (Datalog extension)

Naive approach


- ▶ Numbers in the binary address range → symbols
 - ▶ Numbers outside the range → literals
-
- ▶ False positives: A literal coincides with the binary address range
 - ▶ False negatives: (see paper)

Symbolization: Reducing false positives

- ▶ Collect **additional evidence** (how the number is used) using supporting analyses and heuristics
- ▶ Assign points to candidates:
 - ▷ Symbol
 - ▷ Symbol-Symbol
 - ▷ Strings
 - ▷ Other (data elements with different size)
- ▶ Aggregate points to make a decision

Predicate: `def_used(Adef,Reg,Ause)`

Register `Reg` is defined in `Adef` and used in `Ause`


RDX		40c7f2: <code>mov EDX, 443b08</code>
		40c7f7: <code>xor ECX,ECX</code>
		40c7f9: <code>nop</code>
		40c800: <code>cmp [RDX],RDI</code>
		<code>⋮</code>
		40c808: <code>add RDX,32</code>
		40c80c: <code>cmp RCX,58</code>
		40c810: <code>jle 40c800</code>
		<code>⋮</code>
		443b28: <code>...</code>

Supporting Analyses: Register value analysis

Predicate: $\text{reg_val}(A1, \text{Reg1}, A2, \text{Reg2}, \text{Mult}, \text{Disp})$

$\text{value}(\text{Reg1}@A1) = \text{value}(\text{Reg2}@A2) \times \text{Mult} + \text{Disp}$

```
40c7f2:  mov EDX, 443b08
40c7f7:  xor ECX,ECX
40c7f9:  nop
40c800:  cmp [RDX],RDI
      :      :
40c808:  add RDX,32
40c80c:  cmp RCX,58
40c810:  jle 40c800
      :      :
443b28:  ...
```

 RDX=443b08

Supporting Analyses: Register value analysis

Predicate: $\text{reg_val}(A1, \text{Reg1}, A2, \text{Reg2}, \text{Mult}, \text{Disp})$
 $\text{value}(\text{Reg1}@A1) = \text{value}(\text{Reg2}@A2) \times \text{Mult} + \text{Disp}$

```
40c7f2:  mov EDX, 443b08
40c7f7:  xor ECX,ECX
40c7f9:  nop
40c800:  cmp [RDX],RDI
      :
      :
40c808:  add RDX,32
40c80c:  cmp RCX,58
40c810:  jle 40c800
      :
      :
443b28:  ...
```

$\text{RDX} = ? * 32 + 443b28$

Supporting Analyses: Data access patterns

Predicate: `data_access_pattern(Addr,Size,Mult,Addr2)`

`Addr` is accessed with size `Size` and multiplier `Mult` from `Addr2`

```
40c7f2:  mov EDX, 443b08
```

```
40c7f7:  xor ECX,ECX
```

```
40c7f9:  nop
```

```
40c800:  cmp [RDX],RDI
```

```
⋮
```

```
40c808:  add RDX,32
```

```
40c80c:  cmp RCX,58
```

```
40c810:  jle 40c800
```

```
⋮
```

```
443b28:  ...
```

$RDX = ? * 32 + 443b28$

access Qword with x32 multiplier

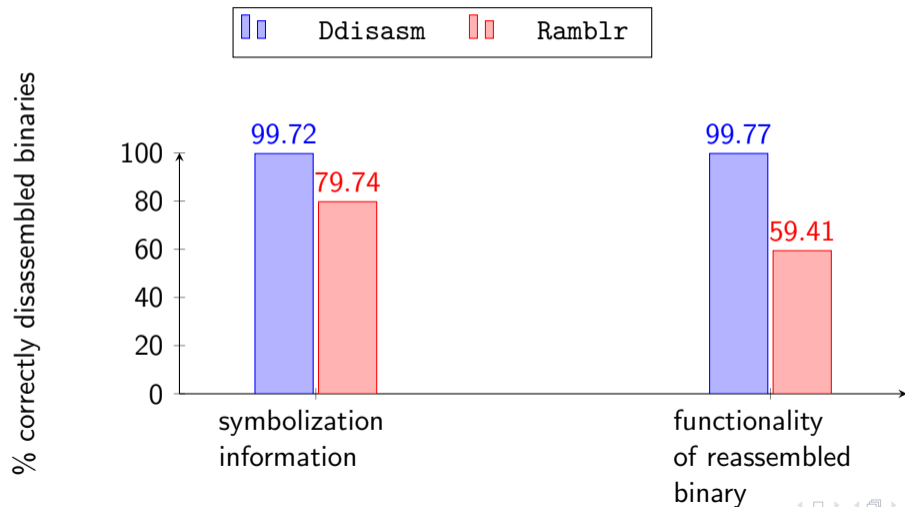
Use supporting analyses to enhance confidence

Candidates in data section

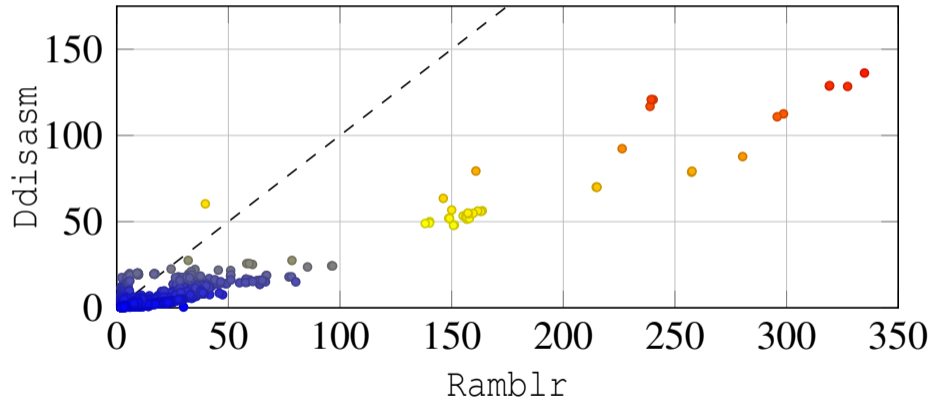
- ▶ Pointer to instruction beginning 👍
- ▶ Data access match 👍
- ▶ Data access conflict 👎
- ▶ ...

- ▶ Our tool `Ddisasm` supports x64 Linux ELF binaries
- ▶ We test our disassembler `Ddisasm` with:
 - ▷ 3 benchmark sets
 - ▷ 7 compiler versions (GCC, Clang and ICC).
 - ▷ 6 compiler optimization flags
- ▶ A total of 7658 binaries
- ▶ Compare to `Ramb1r` state-of-the-art tool in reassembleable disassembly

Results



Disassembly Time in Seconds



- ▶ Reassemblable disassembly is undecidable
- ▶ Practical solutions benefit from analysis and heuristics
- ▶ Datalog works for both:
 - ▷ Express analysis concisely: less error-prone, fast development
 - ▷ Easy to experiment with heuristics expressed as Datalog rules
- ▶ `Ddisasm` is faster and achieves better precision than the state-of-the-art.

Questions?

- ▶ Contact: afloresmontoya@grammatech.com
- ▶ Tool: <https://github.com/GrammaTech/ddisasm>
- ▶ Experimental evaluation: <https://zenodo.org/record/3691736>
- ▶ Rewriting tutorial using Ddisasm :
https://grammatech.github.io/gtirb/md_stack-stamp.html