# Probabilistic Disassembly

Kenneth Adam Miller*, Yonghwi Kwon†, Yi Sun*, Zhuo Zhang*, Xiangyu Zhang*, Zhiqiang Lin‡

\* Department of Computer Science, **Purdue University**, *West Lafayette, USA*

†Department of Computer Science, **University of Virginia**, *Charlottesville, USA*

‡Department of Computer Science and Engineering, **Ohio State University**, *Columbus, USA*

# What is Disassembler?



**Various *semantics lost*:**

(1) Variable names

(2) Data structures (e.g., variable boundaries)
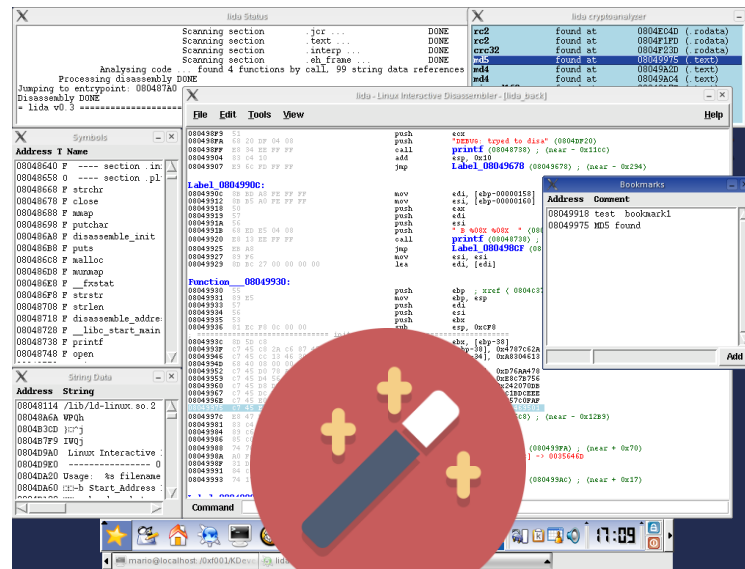
(3) Function names (i.e., indirect call targets)

(4) …

*irreversible* transformation

# What is Disassembler?

**Binary Program** → **Disassembler** → *Human Readable Assembly Code*



*Recover lost semantics with ...*

# Uses of Disassembly

- Reverse-Engineering (e.g., of legacy applications)

- Malware Analysis

- Binary Rewriting for
  - Optimization
  - Security – Software Hardening, Hot Patches, De-bloating (Attack Surface reduction)
  - Instrumentation and Debugging

- …

# Fundamental Problem in Disassembly

When we recover the lost semantics,
***uncertainty*** *is inevitable.*

Especially problematic in recovering
***indirect control flow*** by
*function pointers, virtual tables, switch-case, etc.,*
*and making it difficult to know* **where the code begins**

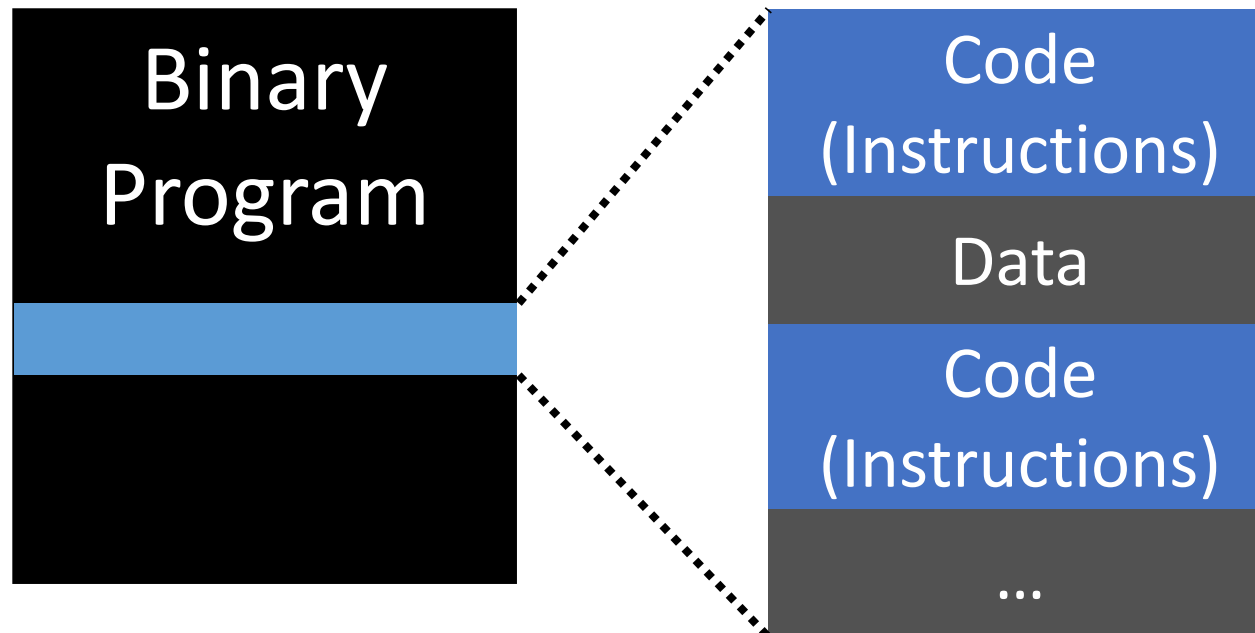# Fundamental Problem in Disassembly

In Von-Neumann Architecture,
*code and data coexist* within the same memory space,
and they could be interleaving

# Fundamental Problem in Disassembly

In Von-Neumann Architecture,
*code and data coexist* within the same memory space,
and they could be interleaving

| Binary Program |
| --- |

Code (Instructions)

Data

Code (Instructions)

...

# Data and Code Coexistence: Linear Sweep

| ADDRESS | RAW BYTE | ASSEMBLY (Ground-truth) |
|---------|----------|--------------------------|
| 0xBBF72 | 66 66 66 66 | DB  66 66 66 66 |
| 0xBBF76 | 66 2E 0F 1F | DB  66 2E 0F 1F |
| 0xBBF7A | 84 00 00 00 | DB  84 00 00 00 |
| 0xBBF7E | 00 00 90 0F | DB  00 00 90 0F |
| 0xBBF82 | 1F 84 00 00 | DB  1F 84 00 00 |
| 0xBBF84 | 00 00 00 00 | DB  00 00 00 00 |
| 0xBBF8A | 00 00 00 00 | DB  00 00 00 00 |
| 0xBBF8E | 00 00 | DB  00 00 |
| 0xBBF90 | 48 8B 05 71 B9 19 00 | mov 0x19B978(rip), rax |
| 0xBBF97 | 41 56 | push r14 |
| ... | ... | ... |
| 0xBBFA2 | 48 8D 50 10 | lea 16(rax), rdx |
| 0xBBFA6 | 48 05 90 00 00 00 | add 144, rax |
| 0xBBFAC | 48 89 47 10 | mov rax, 16(rdi) |
| 0xBBFB0 | 48 89 17 | mov rdx, (rdi) |

| ADDRESS | RAW BYTE | Linear Sweep |
|---------|----------|--------------|
| 0xBBF72 | 66 66 66 66 | nop cs:(rax,rax) |
| 0xBBF76 | 66 2E 0F 1F | ... |
| 0xBBF7A | 84 00 00 00 | ... |
| 0xBBF7C | 00 00 90 0F | ... |
| 0xBBF82 | ... | ... |
| 0xBBF8A | 00 | ... |
| 0xBBF8B | 00 00 | add al, (rax) |
| 0xBBF8D | 00 00 | add al, (rax) |
| 0xBBF8F | 00 48 8B | add cl, -117(rax) |
| 0xBBF92 | 05 71 B9 19 00 | add 1685873, eax |
| 0xBBF97 | 41 56 | push r14 |
| ... | ... | ... |
| 0xBBFA2 | 48 8D 50 10 | lea 16(rax), rdx |
| 0xBBFA6 | 48 05 90 00 00 00 | add 144, rax |
| 0xBBFAC | 48 89 47 10 | mov rax, 16(rdi) |
| 0xBBFB0 | 48 89 17 | mov rdx, (rdi) |

# Data and Code Coexistence: Superset Disassembly

| ADDRESS | RAW BYTE | ASSEMBLY (Ground-truth) |
|---------|----------|--------------------------|
| 0xBBF72 | 66 66 66 66 | DB   66 66 66 66 |
| 0xBBF76 | 66 2E 0F 1F | DB   66 2E 0F 1F |
| 0xBBF7A | 84 00 00 00 | DB   84 00 00 00 |
| 0xBBF7E | 00 00 90 0F | DB   00 00 90 0F |
| 0xBBF82 | 1F 84 00 00 | DB   1F 84 00 00 |
| 0xBBF84 | 00 00 00 00 | DB   00 00 00 00 |
| 0xBBF8A | 00 00 00 00 | DB   00 00 00 00 |
| 0xBBF8E | 00 00 | DB   00 00 |
| 0xBBF90 | 48 8B 05 71 B9 19 00 | mov 0x19B978(rip), rax |
| 0xBBF97 | 41 56 | push r14 |
| ... | ... | ... |
| 0xBBFA2 | 48 8D 50 10 | lea 16(rax), rdx |
| 0xBBFA6 | 48 05 90 00 00 00 | add 144, rax |
| 0xBBFAC | 48 89 47 10 | mov rax, 16(rdi) |
| 0xBBFB0 | 48 89 17 | mov rdx, (rdi) |

| ADDRESS | RAW BYTE | Superset Disassembly |
|---------|----------|----------------------|
| 0xBBF72 | 66 66 .. 00 | nop cs:(rax,rax) |
| 0xBBF73 | 66 66 .. 1F | nop cs:(rax+...) |
| 0xBBF74 | 66 66 .. 84 | nop cs:(rax+...) |
| ... | ... | ... |
| 0xBBF8C | 00 00 | add al, (rax) |
| 0xBBF8D | 00 00 | add al, (rax) |
| 0xBBF8E | 00 00 | add al, (rax) |
| 0xBBF8F | 00 48 8B | add cl, -117(rax) |
| 0xBBF90 | 48 8B 05 71 B9 19 00 | mov 0x19B978(rip), rax |
| 0xBBF91 | 8B 05 .. 00 | mov 0x19b978(rip), eax |
| 0xBFF92 | 05 71 .. 00 | add 0x1685873, eax |
| ... | ... | ... |
| 0xBBF97 | 41 56 | push r14 |
| 0xBBF98 | 56 | push rsi |
| ... | ... | ... |
| 0xBBFB0 | 48 89 47 10 | mov rdx, (rdi) |
| 0xBBFB1 | 48 89 17 | mov edx, (rdi) |
| ... | ... | ... |

# Data and Code Coexistence: Superset Disassembly

| ADDRESS | RAW BYTE | ASSEMBLY (Ground-truth) |
|---------|----------|-------------------------|
| 0xBBF72 | 66 66 66 66 | DB   66 66 66 66 |
| 0xBBF76 | 66 2E 0F 1F | DB   66 2E 0F 1F |
| 0xBBF7A | 84 00 00 00 | DB   84 00 00 00 |
| 0xBBF7 | | |
| 0xBBF8 | | |
| 0xBBF8 | | |
| 0xBBF8 | | |
| 0xBBF8 | | |
| 0xBBF9 | | ax |
| 0xBBF9 | | |
| .. | | |
| 0xBBFA | | |
| 0xBBFA | | |
| | 00 00 | |
| 0xBBFAC | 48 89 47 10 | mov rax, 16(rdi) |
| 0xBBFB0 | 48 89 17 | mov rdx, (rdi) |

*No False Negatives!*

| ADDRESS | RAW BYTE | Superset Disassembly |
|---------|----------|----------------------|
| 0xBBF72 | 66 66 .. 00 | nop cs:(rax,rax) |
| 0xBBF73 | 66 66 .. 1F | nop cs:(rax+...) |
| 0xBBF74 | 66 66 .. 84 | nop cs:(rax+...) |
| ... | ... | ... |
| 0xBBF8C | 00 00 | add al, (rax) |
| 0xBBF8D | 00 00 | add al, (rax) |
| 0xBBF8E | 00 00 | add al, (rax) |
| 0xBBF8F | 00 48 8B | add cl, -117(rax) |
| 0xBBF90 | 48 8B 05 71 B9 19 00 | mov 0x19B978(rip), rax |
| 0xBBF91 | 8B 05 .. 00 | mov 0x19b978(rip), eax |
| 0xBFF92 | 05 71 .. 00 | add 0x1685873, eax |
| ... | ... | ... |
| 0xBBF97 | 41 56 | push r14 |
| 0xBBF98 | 56 | push rsi |
| ... | ... | ... |
| 0xBBFB0 | 48 89 47 10 | mov rdx, (rdi) |
| 0xBBFB1 | 48 89 17 | mov edx, (rdi) |
| ... | ... | ... |

# Data and Code Coexistence: Superset Disassembly

| ADDRESS | RAW BYTE | ASSEMBLY (Ground-truth) |
|---|---|---|
| 0xBBF72 | 66 66 66 66 | DB 66 66 66 66 |
| 0xBBF76 | 66 2E 0F 1F | DB 66 2E 0F 1F |
| 0xBBF7A | 84 00 00 00 | DB 84 00 00 00 |

*Many False Positives*

| ADDRESS | RAW BYTE | Superset Disassembly |
|---|---|---|
| 0xBBF72 | 66 66 .. 00 | **nop** cs:(rax,rax) |
| 0xBBF73 | 66 66 .. 1F | **nop** cs:(rax+...) |
| 0xBBF74 | 66 66 .. 84 | **nop** cs:(rax+...) |
| ... | ... | ... |
| 0xBBF8C | 00 00 | **add** al, (rax) |
| 0xBBF8D | 00 00 | **add** al, (rax) |
| 0xBBF8E | 00 00 | **add** al, (rax) |
| 0xBBF8F | 00 48 8B | **add** cl, -117(rax) |
| 0xBBF90 | 48 8B 05 71 B9 19 00 | **mov** 0x19B978(rip), rax |
| 0xBBF91 | 8B 05 .. 00 | **mov** 0x19b978(rip), eax |
| 0xBFF92 | 05 71 .. 00 | **add** 0x1685873, eax |
| ... | ... | ... |
| 0xBBF97 | 41 56 | **push** r14 |
| 0xBBF98 | 56 | **push** rsi |
| ... | ... | ... |
| 0xBBFB0 | 48 89 47 10 | **mov** rdx, (rdi) |
| 0xBBFB1 | 48 89 17 | **mov** edx, (rdi) |
| ... | ... | ... |

| | | |
|---|---|---|
| | 00 00 | |
| 0xBBFAC | 48 89 47 10 | **mov** rax, 16(rdi) |
| 0xBBFB0 | 48 89 17 | **mov** rdx, (rdi) |

# Data and Code Coexistence: Probabilistic Disassembly

| ADDRESS | RAW BYTE | ASSEMBLY (Ground-truth) |
|---|---|---|
| 0xBBF72 | 66 66 66 66 | DB 66 66 66 66 |
| 0xBBF76 | 66 2E 0F 1F | DB 66 2E 0F 1F |
| 0xBBF7A | 84 00 00 00 | DB 84 00 00 00 |
| 0xBBF7E | 00 00 90 0F | DB 00 00 90 0F |
| 0xBBF82 | 1F 84 00 00 | DB 1F 84 00 00 |
| 0xBBF84 | 00 00 00 00 | DB 00 00 00 00 |
| 0xBBF8A | 00 00 00 00 | DB 00 00 00 00 |
| 0xBBF8E | 00 00 | DB 00 00 |
| 0xBBF90 | 48 8B 05 71 B9 19 00 | mov 0x19B978(rip), rax |
| 0xBBF97 | 41 56 | push r14 |
| ... | ... | ... |
| 0xBBFA2 | 48 8D 50 10 | lea 16(rax), rdx |
| 0xBBFA6 | 48 05 90 00 00 00 | add 144, rax |
| 0xBBFAC | 48 89 47 10 | mov rax, 16(rdi) |
| 0xBBFB0 | 48 89 17 | mov rdx, (rdi) |

| ADDRESS | RAW BYTE | Superset Disassembly |
|---|---|---|
| 0xBBF72 | 66 66 .. 00 | nop cs:(rax,rax) |
| 0xBBF73 | 66 66 .. 1F | nop cs:(rax+...) |
| 0xBBF74 | 66 66 .. 84 | nop cs:(rax+...) |
| ... | ... | ... |
| 0xBBF8C | 00 00 | add al, (rax) |
| 0xBBF8D | 00 00 | add al, (rax) |
| 0xBBF8E | 00 00 | add al, (rax) |
| 0xBBF8F | 00 48 8B | add cl, -117(rax) |
| 0xBBF90 | 48 8B 05 71 B9 19 00 | mov 0x19B978(rip), rax |
| 0xBBF91 | 8B 05 .. 00 | mov 0x19b978(rip), eax |
| 0xBFF92 | 05 71 .. 00 | add 0x1685873, eax |
| ... | ... | ... |
| 0xBBF97 | 41 56 | push r14 |
| 0xBBF98 | 56 | push rsi |
| ... | ... | ... |
| 0xBBFB0 | 48 89 47 10 | mov rdx, (rdi) |
| 0xBBFB1 | 48 89 17 | mov edx, (rdi) |
| ... | ... | ... |

# Data and Code Coexistence: Probabilistic Disassembly

| ADDRESS | RAW BYTE | Superset Disassembly | Probabilistic Disassembly |
|---------|----------|----------------------|---------------------------|
| 0xBBF72 | 66 66 .. 00 | **nop** cs:(rax,rax) | 0.04 |
| 0xBBF73 | 66 66 .. 1F | **nop** cs:(rax+...) | 0.04 |
| 0xBBF74 | 66 66 .. 84 | **nop** cs:(rax+...) | 0.04 |
| ... | ... | ... | ... |
| 0xBBF8C | 00 00 | **add** al, (rax) | **0.695** |
| 0xBBF8D | 00 00 | **add** al, (rax) | 0.04 |
| 0xBBF8E | 00 00 | **add** al, (rax) | **0.695** |
| 0xBBF8F | 00 48 8B | **add** cl, -117(rax) | 0.04 |
| 0xBBF90 | 48 8B 05 71 B9 19 00 | **mov** 0x19B978(rip), rax | **0.695** |
| 0xBBF91 | 8B 05 .. 00 | **mov** 0x19b978(rip), eax | 0.04 |
| 0xBFF92 | 05 71 .. 00 | **add** 0x1685873, eax | 0.04 |
| ... | ... | ... | ... |
| 0xBBF97 | 41 56 | **push** r14 | **0.94** |
| 0xBBF98 | 56 | **push** rsi | 0.06 |
| ... | ... | ... | ... |
| 0xBBFB0 | 48 89 47 10 | **mov** rdx, (rdi) | **1.0 (approx.)** |
| 0xBBFB1 | 48 89 17 | **mov** edx, (rdi) | 0.0 (approx.) |
| ... | ... | ... | ... |

# Data and Code Coexistence: Probabilistic Disassembly

| ADDRESS | RAW BYTE | Superset Disassembly | Probabilistic Disassembly |
|---------|----------|---------------------|---------------------------|
| 0xBBF72 | 66 66 .. 00 | **nop** cs:(rax,rax) | 0.04 |
| 0xBBF73 | 66 66 .. 1F | **nop** cs:(rax+...) | 0.04 |
| 0xBBF74 | 66 66 .. 84 | **nop** cs:(rax+...) | 0.04 |
| ... | ... | ... | ... |
| 0xBBF8C | 00 00 | **add** al, (rax) | **0.695** |
| 0xBBF8D | 00 00 | **add** al, (rax) | 0.04 |
| 0xBBF8E | 00 00 | **add** al, (rax) | **0.695** |
| 0xBBF8F | 00 48 8B | **add** cl, -117(rax) | 0.04 |
| 0xBBF90 | 48 8B 05 71 B9 19 00 | **mov** 0x19B978(rip), rax | **0.695** |
| 0xBBF91 | 8B 05 .. 00 | **mov** 0x19b978(rip), eax | 0.04 |
| 0xBFF92 | 05 71 .. 00 | **add** 0x1685873, eax | 0.04 |
| ... | ... | ... | ... |
| 0xBBF97 | 41 56 | **push** r14 | **0.94** |
| 0xBBF98 | 56 | **push** rsi | 0.06 |
| ... | ... | ... | ... |
| 0xBBFB0 | 48 89 47 10 | **mov** rdx, (rdi) | **1.0 (approx.)** |
| 0xBBFB1 | 48 89 17 | **mov** edx, (rdi) | 0.0 (approx.) |
| ... | ... | ... | ... |

# Data and Code Coexistence: Probabilistic Disassembly

| ADDRESS | RAW BYTE | Superset Disassembly | Probabilistic Disassembly |
|---------|----------|---------------------|---------------------------|
| 0xBBF72 | 66 66 .. 00 | **nop** cs:(rax,rax) | 0.04 |
| 0xBBF73 | 66 66 .. 1F | **nop** cs:(rax+...) | 0.04 |
| 0xBBF74 | 66 66 .. 84 | **nop** cs:(rax+...) | 0.04 |
| ... | ... | ... | ... |
| 0xBBF8C | 00 00 | **add** al, (rax) | **0.695** |
| 0xBBF8D | 00 00 | **add** al, (rax) | 0.04 |
| 0xBBF8E | 00 00 | **add** al, (rax) | **0.695** |
| 0xBBF8F | 00 48 8B | **add** cl, -117(rax) | 0.04 |
| 0xBBF90 | 48 8B 05 71 B9 19 00 | **mov** 0x19B978(rip), rax | **0.695** |
| 0xBBF91 | 8B 05 .. 00 | **mov** 0x19b978(rip), eax | 0.04 |
| 0xBFF92 | 05 71 .. 00 | **add** 0x1685873, eax | 0.04 |
| ... | ... | ... | ... |
| 0xBBF97 | 41 56 | **push** r14 | **0.94** |
| 0xBBF98 | 56 | **push** rsi | 0.06 |
| ... | ... | ... | ... |
| 0xBBFB0 | 48 89 47 10 | **mov** rdx, (rdi) | **1.0 (approx.)** |
| 0xBBFB1 | 48 89 17 | **mov** edx, (rdi) | 0.0 (approx.) |
| ... | ... | ... | ... |

# State-of-the-art: Where we are

Existing disassemblers have
*various limitations*

| Disassembler | False Negatives | False Positives |
|---|---|---|
| Linear sweep | Some | Substantial |
| Traversal | Substantial | None |
| Superset | None | Bloated |
| Shingled Graph | Some | Some |
| BAP ByteWeight | Some | Some |
| *Probabilistic Disassembly* | *None\** | *Some\** |

*Desirable*

*Undesirable*

**\*With probabilistic guarantees**

# Key Ideas: *Hints* of true instructions

*Hint 1:* **Control flow convergence**    *Hint 2*: **Control flow crossing**



**Unlikely** *happen in disassembled* **code from data (e.g., random) bytes**

# Key Ideas: *Hints* of true instructions

## *Hint 3:* Register Def-use relation

(a) True Instructions (Reg. def-use)

| ADDRESS | RAW BYTE | ASSEMBLY |
|---------|----------|----------|
| 0x4005CB | 48 89 C2 | mov rax, rdx |
| 0x4005CE | 48 03 55 F8 | add -0x8(rbp), rdx |
| 0x4005D2 | 8B 45 F4 | mov -0xC(rbp), eax |
| 0x4005D5 | 89 02 | mov eax, (rdx) |

(b) String  (No def-use)

| ADDRESS | RAW BYTE | ASSEMBLY |
|---------|----------|----------|
| 0x400370 | 00 5F 5F | add bl, 0x5F(rdi) |
| 0x400373 | 67 6D | insl cl, (rax) |
| 0x400375 | 6F | outsl ds:(rsi),(dx) |
| 0x400376 | 6E | outsb ds:(rsi),(dx) |

(c) Jump Table  (*Memory* def-use)

| ADDRESS | RAW BYTE | ASSEMBLY |
|---------|----------|----------|
| 0x40040D | 00 00 | add al, (rax) |
| 0x40040F | 00 00 | add cl, (rax) |
| 0x400411 | 10 60 00 | adc ah, 0x0(rax) |

# Computing Probabilistic Constraints

**Reasoning Rules**: probabilistic inference rules

1. If an instruction starting at an address is **valid**, its *successor along control flow* **must be a valid** instruction.

$$C1: \quad x \xrightarrow{1.0} x_{succ}$$

2. If an instruction starting at an address is **valid**, its *preceding instruction* **may** *be a valid instruction.*

$$C2: \quad x \xrightarrow{p} x_{pre}$$

3. If an instruction starting at an address is **valid**, *all the other addresses starting inside its body* are **invalid**

…

*Hints and Reasoning rules* are encoded as *probabilistic constraints.*
Each instruction's probability is propagated to *its control flow successor(s)*

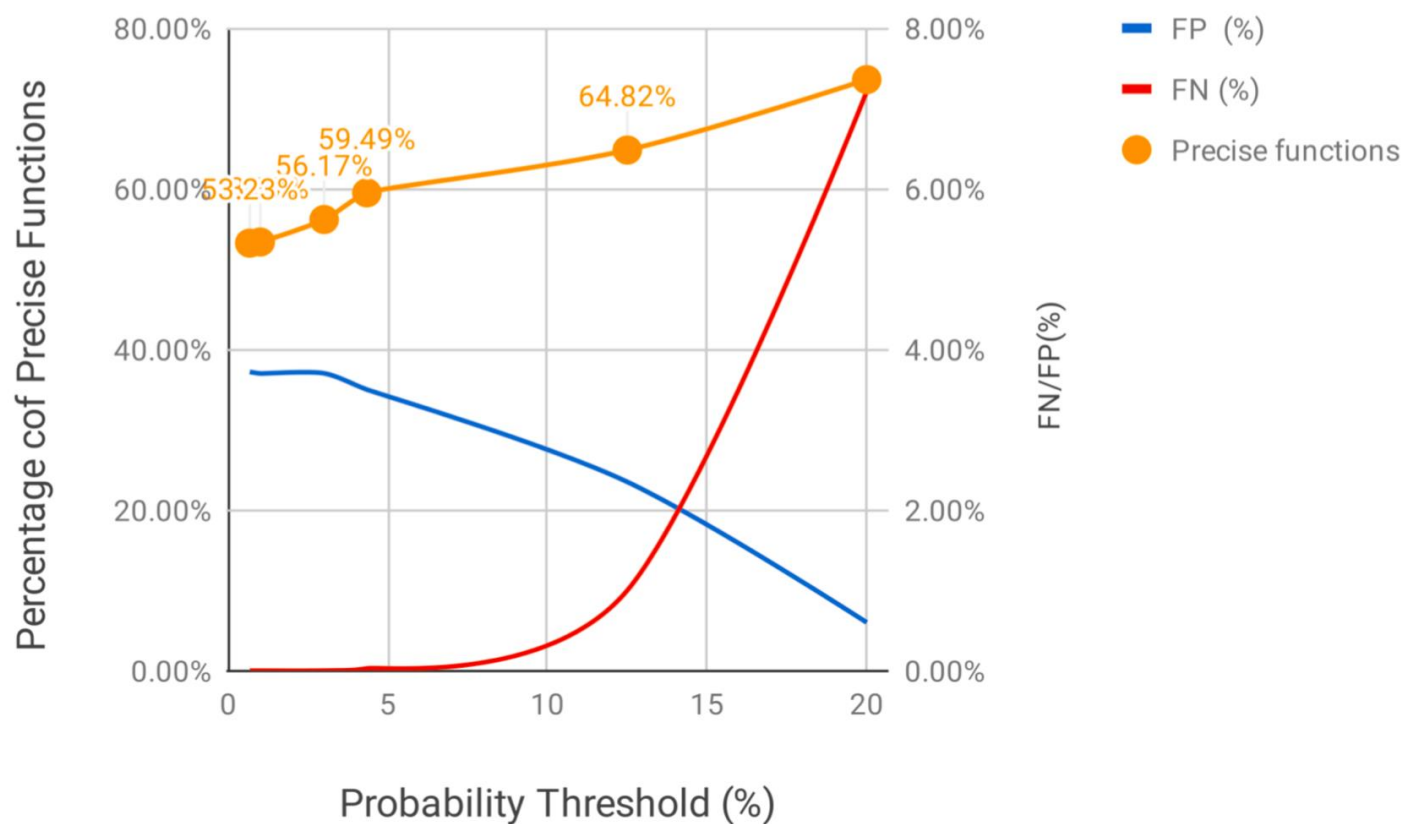# Implementation and Evaluation

- **Implementation**
  - ELF binaries: based on **BAP** (by CMU)
  - PE binaries: based on **Capstone**
  - Released at
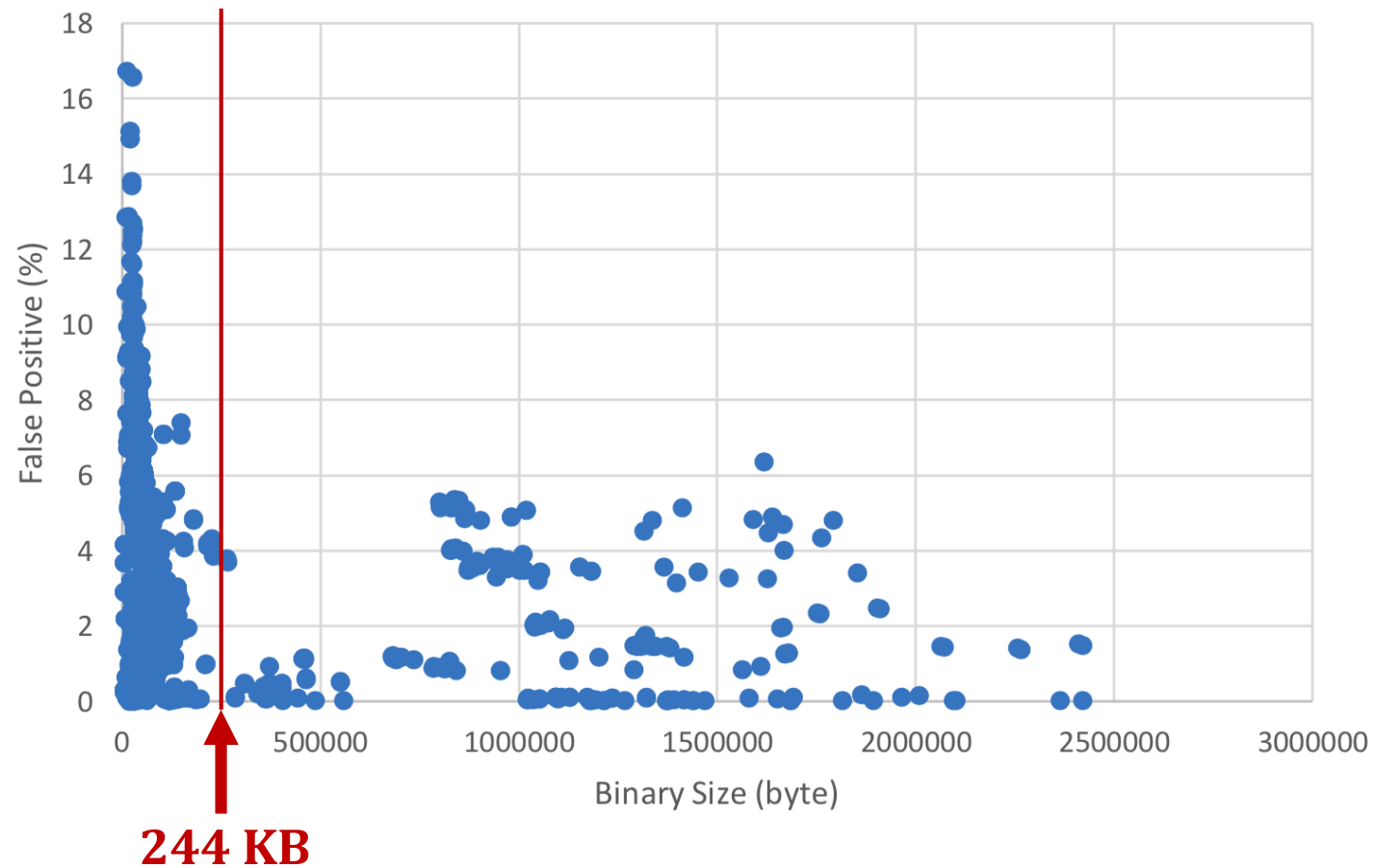    **https://github.com/KennethAdamMiller/superset_disassembler**

- **Evaluation**
  - 2048 ELF binaries
  - SPEC2000 PE binaries

# ELF Results: Tradeoffs of Threshold

# False Positive wrt Binary Size

# Superset Disassembly vs. Probabilistic Disassembly

*Superset Disassembly* + Binary Rewriter* vs.
*Probabilistic Disassembly* + Binary Rewriter*
* The same binary rewriter. No added instructions

| Program | Superset Disassembly | | | Probabilistic Disassembly | | |
| --- | --- | --- | --- | --- | --- | --- |
| | FP | Size (rewritten/orig) | Exec. time (rewritten/orig) | FP | Size (rewritten/orig) | Exec. time (rewritten/orig) |
| 400.perlbench | 85.32% | 780% | 116.71% | 11.29% | 427% | 117.74% |
| 401.bzip2 | 84.65% | 779% | 105.49% | 6.57% | 400% | 97.30% |
| 403.gcc | 88.03% | 751% | 104.60% | 11.33% | 409% | 101.71% |
| 429.mcf | 84.72% | 749% | 104.02% | 4.60% | 399% | 104.74% |
| 445.gobmk | 90.27% | 727% | 103.43% | 6.20% | 372% | 97.30% |
| 456.hmmer | 82.71% | 779% | 99.14% | 6.64% | 411% | 94.12% |
| 458.sjeng | 87.08% | 756% | 98.83% | 7.61% | 407% | 92.76% |
| 462.libquantum | 80.96% | 758% | 100.42% | 4.04% | 400% | 96.94% |
| 464.h264ref | 82.36% | 781% | 100.39% | 2.41% | 395% | 94.57% |
| 471.omnetpp | 85.02% | 768% | 105.24% | 9.82% | 420% | 108.4% |
| 473.astar | 81.46% | 761% | 94.28% | 3.90% | 402% | 93.24% |
| Avg | 84.8% | 763% | 103.0% | 6.8% | 404% | 99.9% |

TABLE II: Superset Disassembly vs Probabilistic Disassembly

# Related Work

- Linear Sweep Disassembly
  - SecondWrite (WCRE'13), D. Andriesse et al. (Security'15), …
- Traversal based Disassembly
  - IDA, BAP, N. E. Rosenblum et al. (AAAI'08), BYTEWEIGHT (Security'14), E. C. R. Shin et al. (Usenix'15), D. Andriesse et al. (Security'15), R. Qiao et al (DSN'17), …
- Superset Disassembly (NDSS'18)
- Probabilistic Inference in Program Analysis
  - G. K. Baah et al. (ISSTA'08), Merlin (PLDI'09), Dimsum (NDSS'12), …
- Machine Learning for Binary Analysis
  - R. Wartell (PKDD'11), Shingled graph disassembly (PAKDD'14), …
- Dynamic Disassembly
  - D. Bruening et al. (CGO'03), H. Patil et al. (MICRO'04), BIRD, …

# Discussions and Questions

- Limitations:
  - Our benchmarks may not represent all possible real-world binaries
  - We focus on compiler generated binaries. Hand-crafted binaries/obfuscated binaries are not tested.
    - However, we believe semantic hints still exist in such code

- Source code available!
  **https://github.com/KennethAdamMiller/superset_disassembler**