# Superset Disassembly: Statically Rewriting x86 Binaries Without Heuristics

Erick Bauman, Zhiqiang Lin, Kevin W. Hamlen

# Motivation

- Static binary rewriting is a core technology for many systems and security applications
- Many static rewriters have been developed over the past decades
- These tools rely on various assumptions and heuristics
- MULTIVERSE: the first heuristic-free static binary rewriter

# Challenge 1:Recognizing and relocating static addresses

```
 1 // gcc -m32 -o sort cmp.o fstring.o sort.c
 2 #include <stdio.h>
 3 #include <unistd.h>
 4
 5 extern char *array[6];
 6 int gt(void *, void *);
 7 int lt(void *, void *);
 8 char* get_fstring(int select);
 9
10 void mode1(void){
11     qsort(array, 5, sizeof(char*), gt);       C4
12 }
13 void mode2(void){
14     qsort(array, 5, sizeof(char*), lt);       C4
15 }
16                                               C1
17 void (*modes[2])() = {mode1, mode2};
18
19 void main(void) {
20     int p = getpid() & 1;
21     printf(get_fstring(0),p);
22     (*modes[p])();                            C2
23     print_array();
24 }
```

# Challenge 1:Recognizing and relocating static addresses

```
Hex dump of section '.data':
  0x0804a01c 00000000 00000000 70870408 74870408 ........p...t...
  0x0804a02c 78870408 7d870408 81870408 00000000 x...}...........
  0x0804a03c f4850408 20860408                    .... ...        C1
```

(f) Hexdump of .data section

# Keeping original data space intact

- No need to modify data addresses if data unchanged

- Keep read-only copy of code for inline data in original code section

# Challenge 2: Handling dynamically computed memory addresses

```
 1 // gcc -m32 -o sort cmp.o fstring.o sort.c
 2 #include <stdio.h>
 3 #include <unistd.h>
 4
 5 extern char *array[6];
 6 int gt(void *, void *);
 7 int lt(void *, void *);
 8 char* get_fstring(int select);
 9
10 void mode1(void){
11     qsort(array, 5, sizeof(char*), gt);                          C4
12 }
13 void mode2(void){
14     qsort(array, 5, sizeof(char*), lt);                          C4
15 }
16                                                                  C1
17 void (*modes[2])() = {mode1, mode2};
18
19 void main(void){
20     int p = getpid() & 1;
21     printf(get_fstring(0),p);
22     (*modes[p])();                                               C2
23     print_array();
24 }
```

# Challenge 2: Handling dynamically computed memory addresses

```
804864c <main>:
...
8048678:    e8 73 fd ff ff          call    80483f0 <printf@plt>
804867d:    8b 44 24 1c             mov     0x1c(%esp),%eax
8048681:    8b 04 85 3c a0 04 08    mov     0x804a03c(,%eax,4),%eax
8048688:    ff d0                   call    *%eax
...
```

C2

(d) Partial binary code of sort

# Creating mapping from old code space to rewritten code space

- Do not attempt to identify original addresses to rewrite

- Ignore how address is computed; only focus on final target

- Rewrite all iCFTs to use mapping to dynamically translate address on use

# Challenge 3: Differentiating code and data

```
 1 ;nasm -f elf fstring.asm
 2 BITS 32
 3 GLOBAL get_fstring
 4 SECTION .text
 5 get_fstring:
 6     mov eax,[esp+4]
 7     cmp eax,0
 8     jz after
 9     mov eax,msg2
10     ret
11 msg1:
12     db 'mode: %d', 10, 0          C3
13 msg2:
14     db '%s', 10, 0                C3
15 after:
16     mov eax,msg1
17     ret
```

(b) Source code of `fstring.asm`

# Challenge 3: Differentiating code and data



(d) Partial binary code of sort

# Brute force disassembling of all possible code

- Disassemble every offset
- All intended code will be within resulting superset

# Challenge 4: Handling function pointer arguments

```
 1 // gcc -m32 -o sort cmp.o fstring.o sort.c
 2 #include <stdio.h>
 3 #include <unistd.h>
 4
 5 extern char *array[6];
 6 int gt(void *, void *);
 7 int lt(void *, void *);
 8 char* get_fstring(int select);
 9
10 void mode1(void){
11     qsort(array, 5, sizeof(char*), gt);       C4
12 }
13 void mode2(void){
14     qsort(array, 5, sizeof(char*), lt);       C4
15 }
16                                               C1
17 void (*modes[2])() = {mode1, mode2};
18
19 void main(void){
20     int p = getpid() & 1;
21     printf(get_fstring(0),p);
22     (*modes[p])();                            C2
23     print_array();
24 }
```

# Challenge 4: Handling function pointer arguments



```
...
80485a0 <gt>:
80485a0:    53                              push    %ebx
...
80485f4 <mode1>:
...
80485fa:    c7 44 24 0c a0 85 04    movl    $0x80485a0,0xc(%esp)
8048601:    08
8048602:    c7 44 24 08 04 00 00    movl    $0x4,0x8(%esp)
8048609:    00
804860a:    c7 44 24 04 05 00 00    movl    $0x5,0x4(%esp)
8048611:    00
8048612:    c7 04 24 24 a0 04 08    movl    $0x804a024,(%esp)
8048619:    e8 12 fe ff ff          call    8048430 <qsort@plt>
...
```

(d) Partial binary code of sort

# Rewriting all user level code including libraries

- Hard to automatically identify all function pointer arguments

- Instead, rewrite everything

- Use mapping(from solution 2) to translate callback upon use

# Challenge 5: Handling PIC

```
 1 // gcc -m32 -c -o cmp.o cmp.c -fPIC -O2
 2 #include <stdio.h>
 3 #include <stdlib.h>
 4 #include <string.h>
 5
 6 char *array[6] = {"foo", "bar", "quuz", "baz", "flux"};        C1
 7 char* get_fstring(int select);
 8
 9 void print_array() {
10     int i;
11     for (i = 0; i < 5; i++) {
12         fprintf(stdout, get_fstring(1), array[i]);             C5
13     }
14 }
15 int lt(void *a, void *b) {
16     return strcmp(*(char **) a, *(char **)b);
17 }
18
19 int gt(void *a, void *b) {
20     return strcmp(*(char **) b, *(char **)a);
21 }
```

(c) Source code of `cmp.c`

# Challenge 5: Handling PIC



```
8048510 <print_array>:
...
8048515:    53                          push    %ebx
8048516:    e8 b1 00 00 00              call    80485cc <__i686.get_pc_thunk.bx>
804851b:    81 c3 d9 1a 00 00           add     $0x1ad9,%ebx
8048521:    83 ec 1c                    sub     $0x1c,%esp
8048524:    8b ab fc ff ff ff           mov     -0x4(%ebx),%ebp
...
80485a0 <gt>:
80485a0:    53                          push    %ebx
...
80485cc <__i686.get_pc_thunk.bx>:
80485cc:    8b 1c 24                    mov     (%esp),%ebx
80485cf:    c3                          ret
```

(d) Partial binary code of **sort**
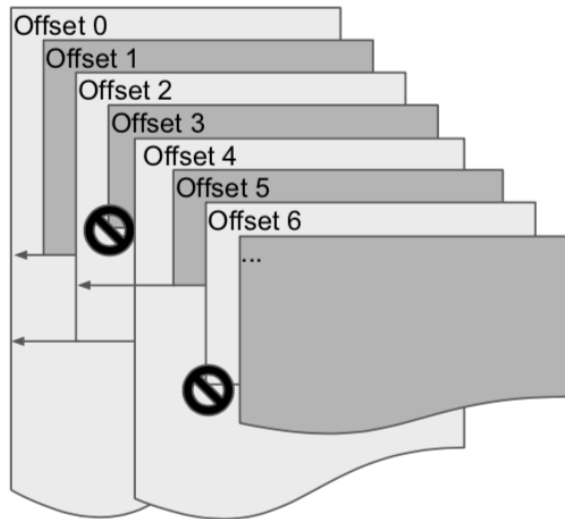
# Rewriting all call instruction

- For x86-32 instructions, only call reveals instruction pointer

- Rewrite call to push/jmp and push old return address

- Offsets computed based on old address

- From solution 2, rewritten ret instructions translate return address with mapping

# Superset Disassembler

**Algorithm 1:** Superset Disassembly

**input** : empty two-dimensional list *instructions*
**input** : string of raw bytes of text section *bytes*
**output** : all disassembled instructions are in *instructions*

1  **for** *start_offset ← 0* **to** *length(bytes)* **do**
2      offset ← start_offset;
3      **while** *legal(offset) and offset ∉ instructions and offset < length(bytes)* **do**
4         instruction ← disassemble(offset);
5         instructions[start_offset][offset] ← instruction;
6         offset ← offset + length(instruction);
7      **if** *offset ∈ instructions* **then**
8         instructions[start_offset][offset] ← "jmp offset";



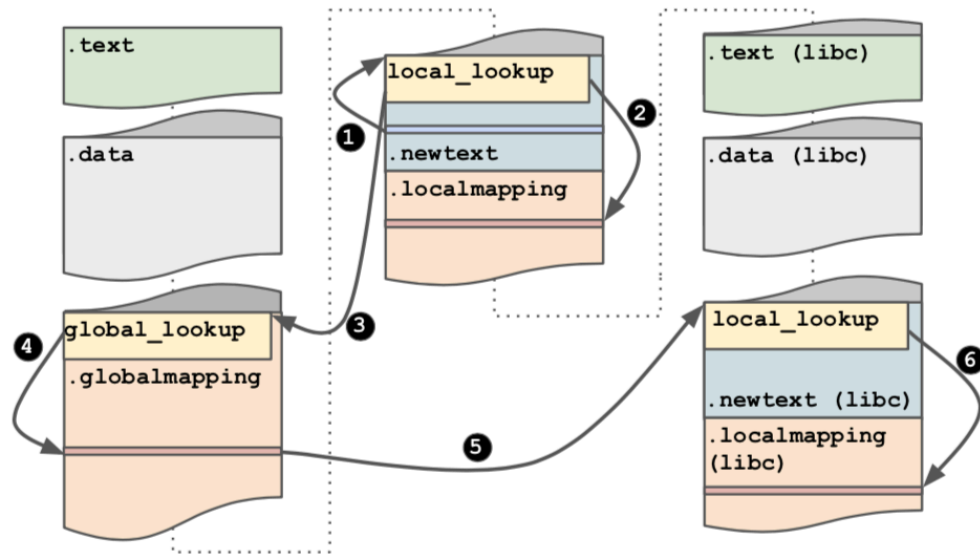: An illustration of our disassembly strategy.

# Mapping



Fig. 4: A mapping lookup example for a rewritten binary dynamically linked with our rewritten `libc`.

# Rewriting

- jcc/jmp/call
- Indirect call

- **jmp:** If the instruction is `jmp [target]`, we rewrite it to the following six instructions:
  ```
  mov [esp-32], eax
  mov eax, target
  call lookup
  mov [esp-4], eax
  mov eax, [esp-32]
  jmp [esp-4]
  ```

# Evaluation

| Benchmark | Dir. Calls | Dir. Jumps | Ind. Calls | Ind. Jumps | Cond. Jumps | Rets | .text (KB) | .newtext (KB) | Size Inc. ($\times$) |
|---|---|---|---|---|---|---|---|---|---|
| 400.perlbench | 30888 | 24778 | 3896 | 4442 | 126876 | 22306 | 1047 | 5146 | 12.88 |
| 401.bzip2 | 1100 | 1050 | 170 | 152 | 7342 | 874 | 55 | 268 | 70.71 |
| 403.gcc | 110122 | 64532 | 8916 | 15680 | 380920 | 45410 | 3225 | 15290 | 10.32 |
| 429.mcf | 276 | 216 | 44 | 78 | 1300 | 250 | 12 | 57 | 202.98 |
| 445.gobmk | 23548 | 14946 | 3550 | 3480 | 117378 | 20918 | 1488 | 6520 | 5.39 |
| 456.hmmer | 8020 | 4942 | 556 | 666 | 28924 | 4106 | 277 | 1279 | 22.56 |
| 458.sjeng | 2566 | 2338 | 256 | 658 | 12236 | 1570 | 132 | 604 | 36.17 |
| 462.libquantum | 1094 | 758 | 94 | 146 | 3376 | 812 | 40 | 181 | 93.73 |
| 464.h264ref | 7124 | 6518 | 1782 | 2000 | 47850 | 6318 | 520 | 2441 | 16.23 |
| 471.omnetpp | 33578 | 10032 | 3830 | 1782 | 51642 | 14326 | 635 | 3029 | 13.49 |
| 473.astar | 912 | 552 | 162 | 160 | 3314 | 750 | 39 | 184 | 92.52 |
| 483.xalancbmk | 115154 | 58678 | 39392 | 14630 | 307122 | 75674 | 3850 | 17369 | 7.60 |
| libc.so.6 | 32798 | 33370 | 9816 | 9012 | 189384 | 32458 | 1735 | 8435 | 9.77 |
| libgcc_s.so.1 | 2158 | 2514 | 374 | 484 | 12862 | 1740 | 112 | 538 | 9.70 |
| libm.so.6 | 5450 | 8870 | 874 | 892 | 21796 | 7406 | 277 | 1268 | 9.51 |
| libstdc++.so.6 | 22456 | 10418 | 4300 | 4008 | 144516 | 15784 | 900 | 4258 | 9.53 |

TABLE I: Statistics of MULTIVERSE rewritten binaries and libraries