

OSPREY: Recovery of Variable and Data Structure via Probabilistic Analysis for Stripped Binary

Zhuo Zhang, Yapeng Ye, Wei You, Guanhong Tao, Wen-chuan Lee,
Yonghwi Kwon, Yousra Aafer, Xiangyu Zhang



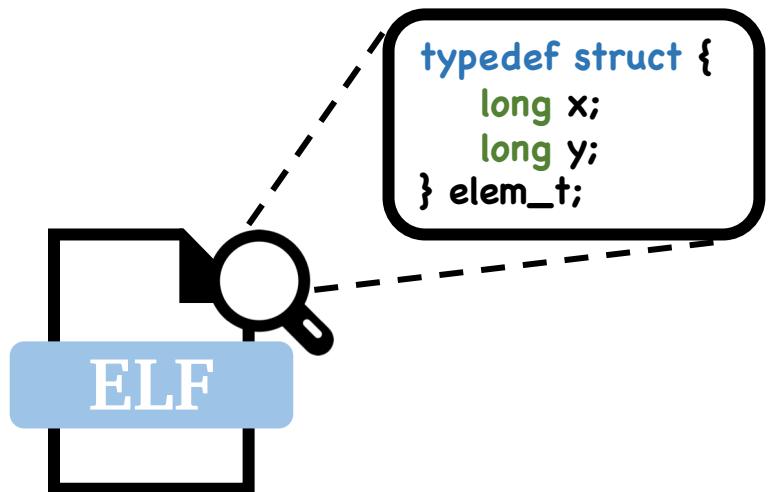
中國人民大學
RENMIN UNIVERSITY OF CHINA



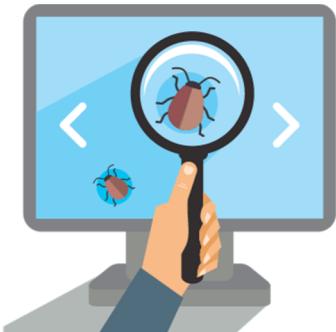
UNIVERSITY
of VIRGINIA



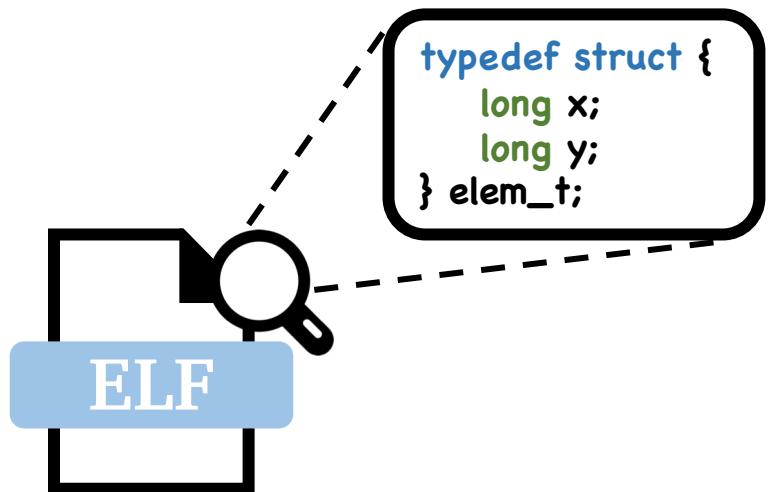
UNIVERSITY OF
WATERLOO



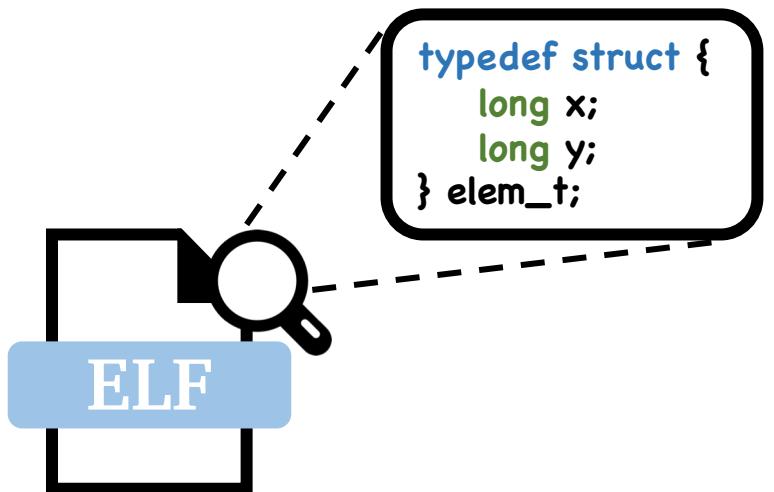
Variable and Data Structure Recovery



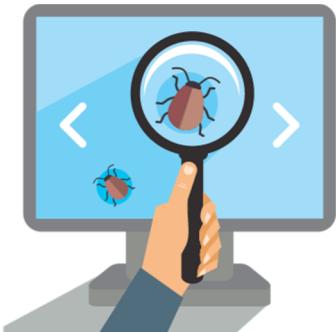
Bug Detection



Variable and Data Structure Recovery



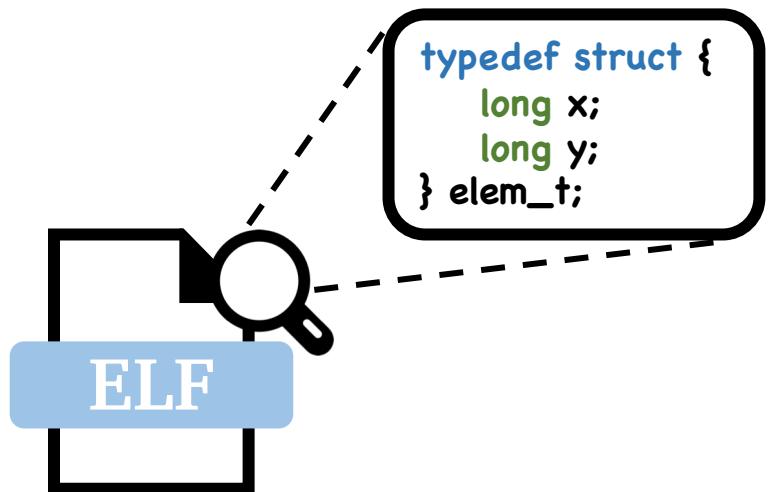
Variable and Data Structure Recovery



Bug Detection



Malware Analysis



Variable and Data Structure Recovery



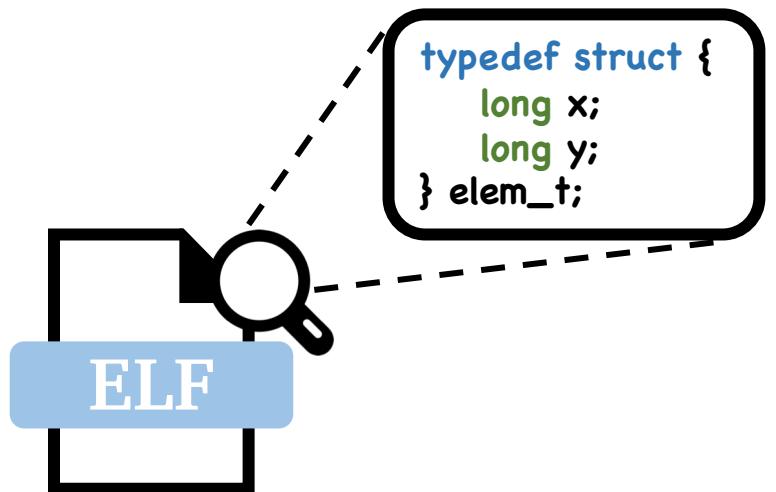
Bug Detection

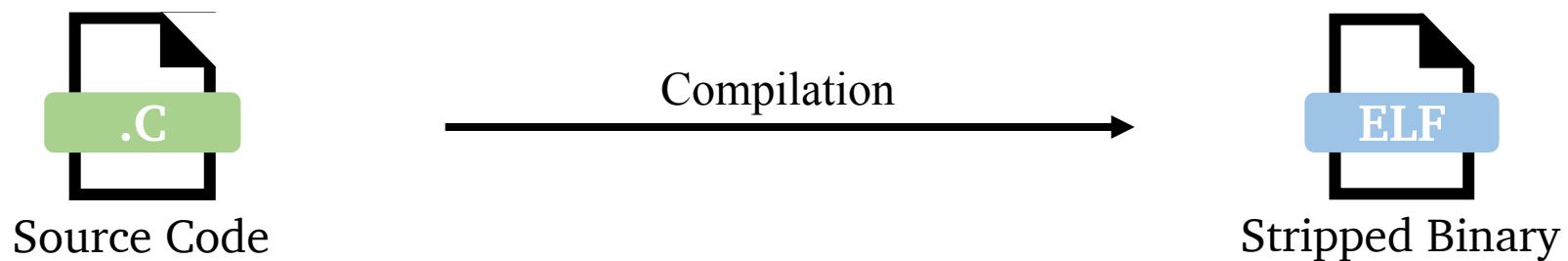


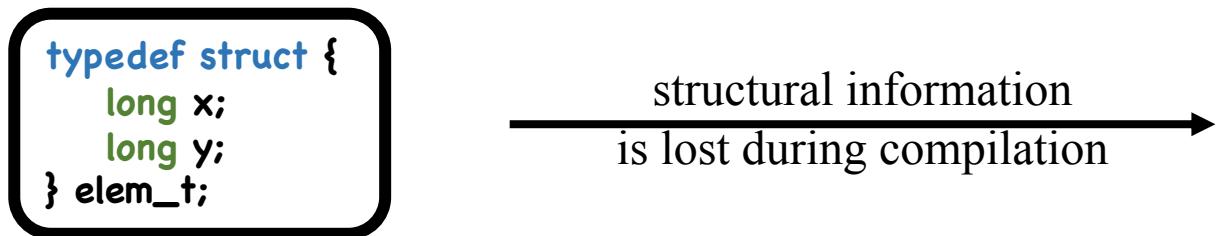
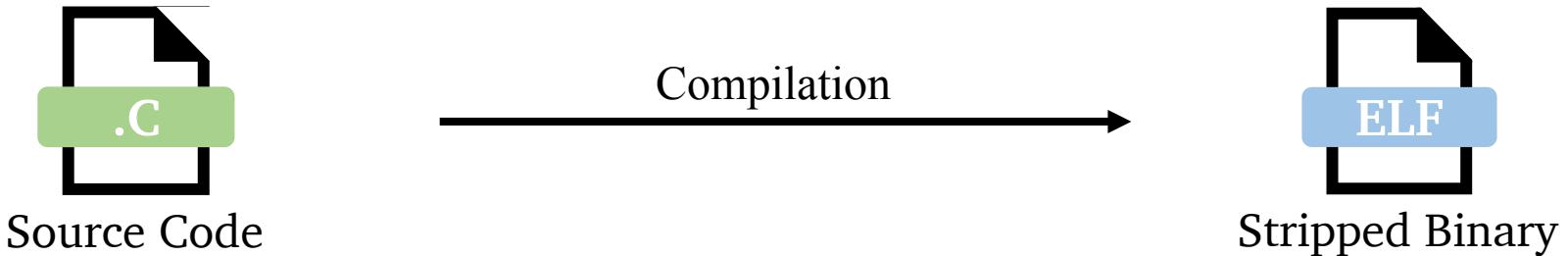
Malware Analysis

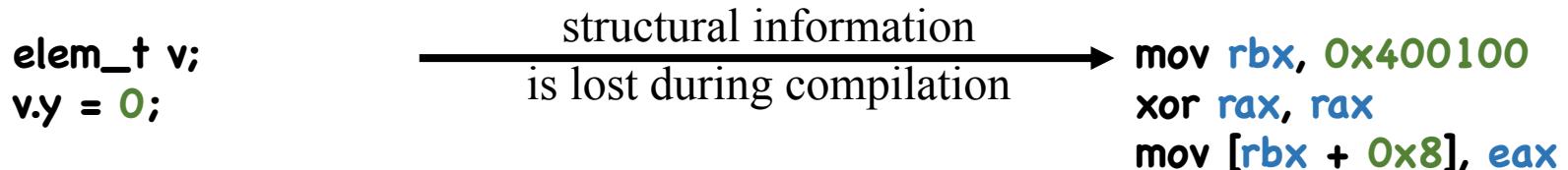
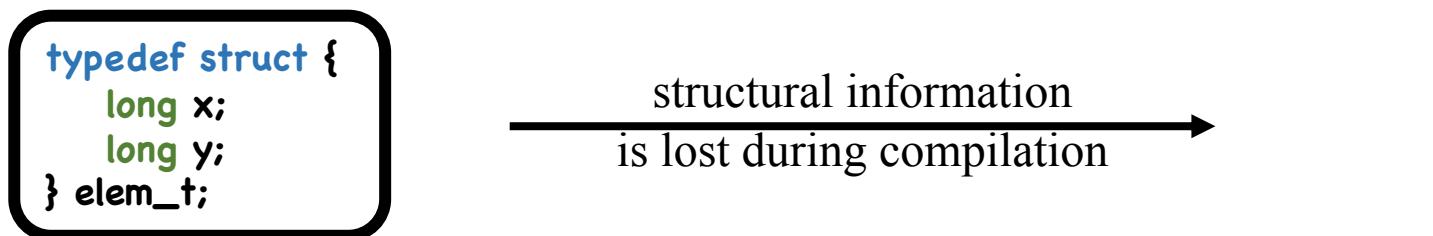
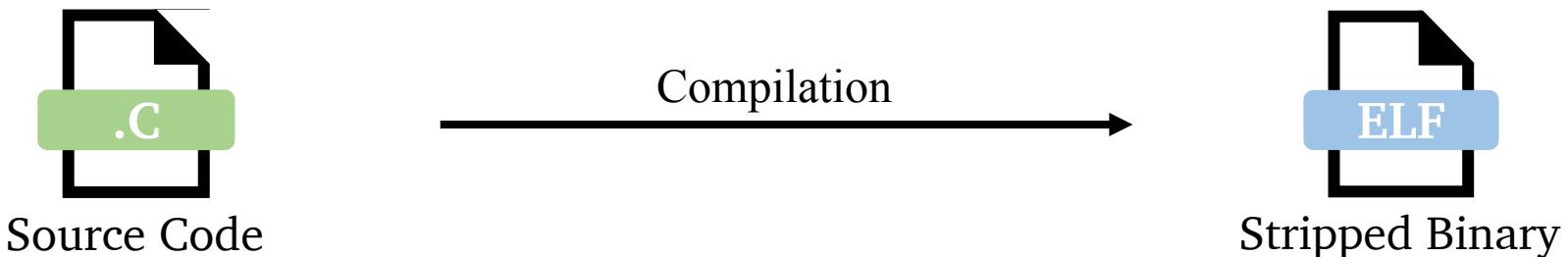


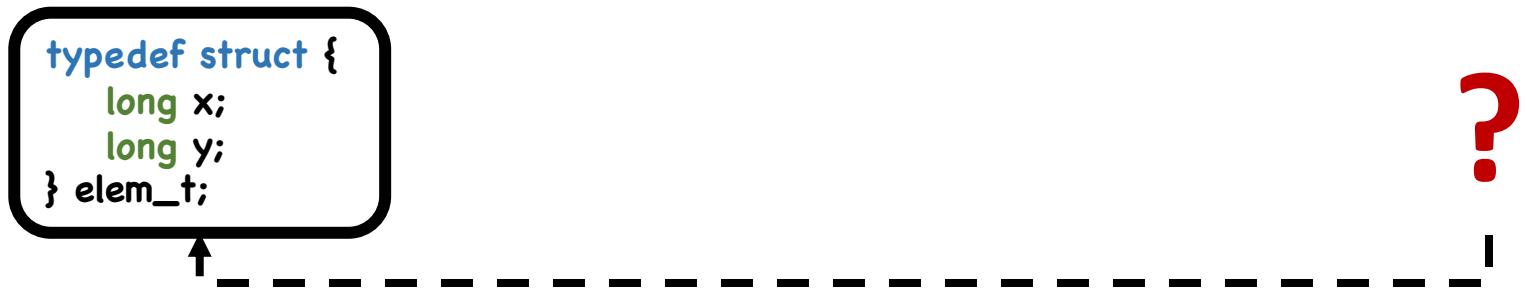
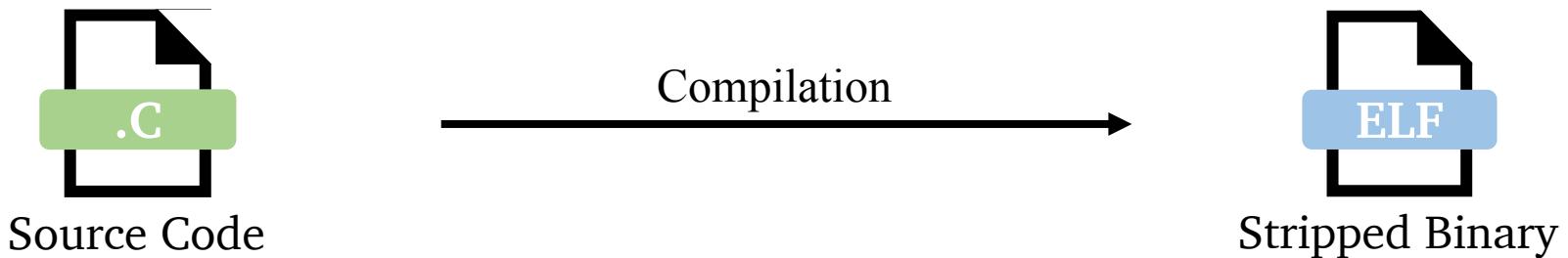
Automatic Exploit Generation











This diagram shows the recovery of assembly code. A dashed line with arrows at both ends connects two horizontal sections. The top section contains assembly code:

```
mov rbx, 0x400100
xor rax, rax
mov [rbx + 0x8], eax
```

The bottom section also contains assembly code, which is identical to the top section:

```
mov rbx, 0x400100
xor rax, rax
mov [rbx + 0x8], eax
```

Motivation Example

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (...)
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Motivation Example

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Motivation Example

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Motivation Example

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (...) {
14.         p = &v;
15.     } else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Motivation Example

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (...)
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Motivation Example

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (...)
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Motivation Example

Variable and data structure recovery is to recover high-level semantic information from the compiled executables.

- Array
- Structure
- Pointer

In the motivation example, we are trying to recover:

- the structure of `elem_t`
- `p` is a pointer to `elem_t`
- `v` is a an `elem_t` located on the stack

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Motivation Example

Variable and data structure recovery is to recover high-level semantic information from the compiled executables.

- Array
- Structure
- Pointer

In the motivation example, we are trying to recover:

- the structure of `elem_t`
- `p` is a pointer to `elem_t`
- `v` is a an `elem_t` located on the stack

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Motivation Example

Variable and data structure recovery is to recover high-level semantic information from the compiled executables.

- Array
- Structure
- Pointer

In the motivation example, we are trying to recover:

- the structure of `elem_t`
- `p` is a pointer to `elem_t`
- `v` is a an `elem_t` located on the stack

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Motivation Example

Variable and data structure recovery is to recover high-level semantic information from the compiled executables.

- Array
- Structure
- Pointer

In the motivation example, we are trying to recover:

- the structure of `elem_t`
- `p` is a pointer to `elem_t`
- `v` is a an `elem_t` located on the stack

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

State-of-the-art: Ghidra and IDA Pro

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

State-of-the-art: Ghidra and IDA Pro

Ghidra and IDA Pro are based on specific code patterns.

- Unreliable due to uncertainty, leading to contradicting results
- Cannot handle complex structures

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

State-of-the-art: Ghidra and IDA Pro

Ghidra and IDA Pro are based on specific code patterns.

- Unreliable due to uncertainty, leading to contradicting results
- Cannot handle complex structures

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

State-of-the-art: Ghidra and IDA Pro

Ghidra and IDA Pro are based on specific code patterns.

- Unreliable due to uncertainty, leading to contradicting results
- Cannot handle complex structures

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

State-of-the-art: Ghidra and IDA Pro

Ghidra and IDA Pro are based on specific code patterns.

- Unreliable due to uncertainty, leading to contradicting results
- Cannot handle complex structures

For example, the **movdqa** and **movups** instruction pair denotes a 128-bit packed floating-point value movement. Ghidra and IDA Pro recognize this and guess the structure is a float-point variable.

```
call malloc
movdqa xmm0, [rsp + 0x8]
movups [rax], xmm0
```



```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (...)

14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

State-of-the-art: Ghidra and IDA Pro

Ghidra and IDA Pro are based on specific code patterns.

- Unreliable due to uncertainty, leading to contradicting results
- Cannot handle complex structures

For example, the **movdqa** and **movups** instruction pair denotes a 128-bit packed floating-point value movement. Ghidra and IDA Pro recognize this and guess the structure is a float-point variable.

```
call malloc
movdqa xmm0, [rsp + 0x8]
movups [rax], xmm0
```



```

01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (...) {
14.         p = &v;
15.     } else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

State-of-the-art: Ghidra and IDA Pro

Ghidra and IDA Pro are based on specific code patterns.

- Unreliable due to uncertainty, leading to contradicting results
- Cannot handle complex structures

For example, the **movdqa** and **movups** instruction pair denotes a 128-bit packed floating-point value movement. Ghidra and IDA Pro recognize this and guess the structure is a float-point variable.

```
typedef union {
    int64 u_0[2];
    int128 u_1;
} elem_t;
```

IDA Pro

```
typedef struct {
    int32 s_0[4];
} elem_t;
```

Ghidra

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (...)

14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

State-of-the-art: Howard

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (...)
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

State-of-the-art: Howard

Howard is a state-of-the-art dynamic analysis tool.

- Use dynamic analysis to collect program behaviors
- Leverage heuristics such as field accesses are performed by first loading the base address of the data structure and then offsetting.

However,

- Its effectiveness hinges on the availability of high quality inputs.
- The heuristics do not always hold.

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

State-of-the-art: Howard

Howard is a state-of-the-art dynamic analysis tool.

- Use dynamic analysis to collect program behaviors
- Leverage heuristics such as field accesses are performed by first loading the base address of the data structure and then offsetting.

However,

- Its effectiveness hinges on the availability of high quality inputs.
- The heuristics do not always hold.

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

State-of-the-art: Howard

Howard is a state-of-the-art dynamic analysis tool.

- Use dynamic analysis to collect program behaviors
- Leverage heuristics such as field accesses are performed by first loading the base address of the data structure and then offsetting.

However,

- Its effectiveness hinges on the availability of high quality inputs.
- The heuristics do not always hold.

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

State-of-the-art: Howard

Howard is a state-of-the-art dynamic analysis tool.

- Use dynamic analysis to collect program behaviors
- Leverage heuristics such as field accesses are performed by first loading the base address of the data structure and then offsetting.

However,

- Its effectiveness hinges on the availability of high quality inputs.
- The heuristics do not always hold.

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

State-of-the-art: Howard

Howard is a state-of-the-art dynamic analysis tool.

- Use dynamic analysis to collect program behaviors
- Leverage heuristics such as field accesses are performed by first loading the base address of the data structure and then offsetting.

However,

- Its effectiveness hinges on the availability of high quality inputs.
- The heuristics do not always hold.

```

01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (...)
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

State-of-the-art: Howard

Howard is a state-of-the-art dynamic analysis tool.

- Use dynamic analysis to collect program behaviors
- Leverage heuristics such as field accesses are performed by first loading the base address of the data structure and then offsetting.

However,

- Its effectiveness hinges on the availability of high quality inputs.
- The heuristics do not always hold.

```

01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (...)

14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

State-of-the-art: Howard

Howard is a state-of-the-art dynamic analysis tool.

- Use dynamic analysis to collect program behaviors
- Leverage heuristics such as field accesses are performed by first loading the base address of the data structure and then offsetting.

However,

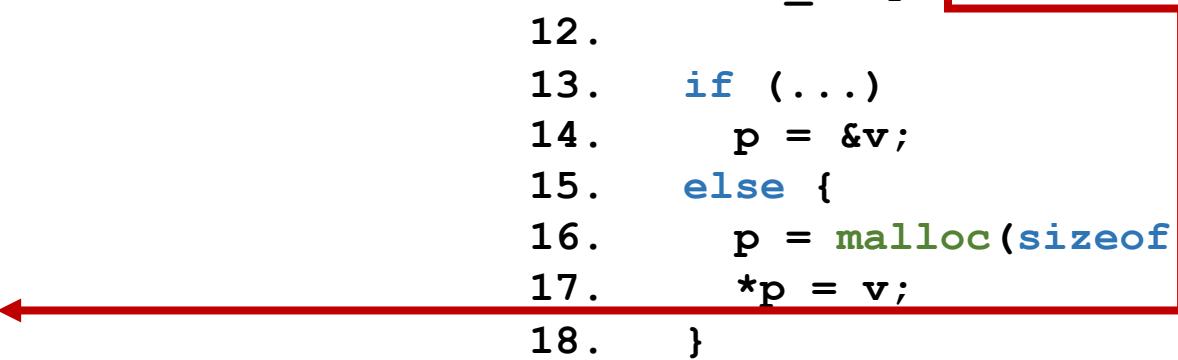
- Its effectiveness hinges on the availability of high quality inputs.
- The heuristics do not always hold.

`mov [rsp + 0x8], 0x0
mov [rsp + 0x10], 0x1`

```

01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (...)

14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```



State-of-the-art: Howard

Howard is a state-of-the-art dynamic analysis tool.

- Use dynamic analysis to collect program behaviors
- Leverage heuristics such as field accesses are performed by first loading the base address of the data structure and then offsetting.

However,

- Its effectiveness hinges on the availability of high quality inputs.
- The heuristics do not always hold.

```
typedef struct {
    int64 s_0;
    int64 s_1;
} elem_t;

elem_t *p;
int64 v0;
int64 v1;
```

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (...)

14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

State-of-the-art: Howard

Howard is a state-of-the-art dynamic analysis tool.

- Use dynamic analysis to collect program behaviors
- Leverage heuristics such as field accesses are performed by first loading the base address of the data structure and then offsetting.

However,

- Its effectiveness hinges on the availability of high quality inputs.
- The heuristics do not always hold.

```
typedef struct {
    int64 s_0;
    int64 s_1;
} elem_t;

elem_t *p;
int64 v0;
int64 v1;
```

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.

06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.

10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.

13.     if (...)

14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.

20.     output(p->x, p->y);
21. }
```

Our Techniques

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (...)
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Our Techniques

Observation 1: while existing techniques mostly focus on memory access patterns (i.e., base addresses and offset values) to identify structures, there are many other program behaviors that can serve as hints to recover data structures.

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Our Techniques

Observation 1: while existing techniques mostly focus on memory access patterns (i.e., base addresses and offset values) to identify structures, there are many other program behaviors that can serve as hints to recover data structures.

Allocation Hint: a heap-allocated variable is *likely* to be a structure or an array of structure.

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Our Techniques

Observation 1: while existing techniques mostly focus on memory access patterns (i.e., base addresses and offset values) to identify structures, there are many other program behaviors that can serve as hints to recover data structures.

Allocation Hint: a heap-allocated variable is *likely* to be a structure or an array of structure.

*p is likely to be structure.

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Our Techniques

Observation 1: while existing techniques mostly focus on memory access patterns (i.e., base addresses and offset values) to identify structures, there are many other program behaviors that can serve as hints to recover data structures.

Allocation Hint: a heap-allocated variable is *likely* to be a structure or an array of structure.

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Our Techniques

Observation 1: while existing techniques mostly focus on memory access patterns (i.e., base addresses and offset values) to identify structures, there are many other program behaviors that can serve as hints to recover data structures.

Allocation Hint: a heap-allocated variable is *likely* to be a structure or an array of structure.

Dataflow Hint: two memory regions connected by a direct dataflow *may* be of the same data structure.

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Our Techniques

Observation 1: while existing techniques mostly focus on memory access patterns (i.e., base addresses and offset values) to identify structures, there are many other program behaviors that can serve as hints to recover data structures.

Allocation Hint: a heap-allocated variable is *likely* to be a structure or an array of structure.

Dataflow Hint: two memory regions connected by a direct dataflow *may* be of the same data structure.

```
p->x = v.x;
p->y = v.y;
```

***p** and **v** may be of the same data structure.

```
01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (...)

14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Our Techniques

Observation 1: while existing techniques mostly focus on memory access patterns (i.e., base addresses and offset values) to identify structures, there are many other program behaviors that can serve as hints to recover data structures.

Allocation Hint: a heap-allocated variable is *likely* to be a structure or an array of structure.

Dataflow Hint: two memory regions connected by a direct dataflow *may* be of the same data structure.

```

01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Our Techniques

Observation 1: while existing techniques mostly focus on memory access patterns (i.e., base addresses and offset values) to identify structures, there are many other program behaviors that can serve as hints to recover data structures.

Allocation Hint: a heap-allocated variable is *likely* to be a structure or an array of structure.

Dataflow Hint: two memory regions connected by a direct dataflow *may* be of the same data structure.

Point-to Hint: if a pointer may point to two memory regions, these regions *may* be of the same data structure.

```

01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Our Techniques

Observation 1: while existing techniques mostly focus on memory access patterns (i.e., base addresses and offset values) to identify structures, there are many other program behaviors that can serve as hints to recover data structures.

Allocation Hint: a heap-allocated variable is *likely* to be a structure or an array of structure.

Dataflow Hint: two memory regions connected by a direct dataflow *may* be of the same data structure.

Point-to Hint: if a pointer may point to two memory regions, these regions *may* be of the same data structure.

p may point to both stack inlined variable **v** and heap-allocated variable **malloc**(...), suggesting they may be of the same structure.

```

01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (...)

14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Our Techniques

Observation 1: while existing techniques mostly focus on memory access patterns (i.e., base addresses and offset values) to identify structures, there are many other program behaviors that can serve as hints to recover data structures.

Allocation Hint: a heap-allocated variable is *likely* to be a structure or an array of structure.

Dataflow Hint: two memory regions connected by a direct dataflow *may* be of the same data structure.

Point-to Hint: if a pointer may point to two memory regions, these regions *may* be of the same data structure.

```

01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Our Techniques

Observation 1: while existing techniques mostly focus on memory access patterns (i.e., base addresses and offset values) to identify structures, there are many other program behaviors that can serve as hints to recover data structures.

Allocation Hint: a heap-allocated variable is *likely* to be a structure or an array of structure.

Dataflow Hint: two memory regions connected by a direct dataflow *may* be of the same data structure.

Point-to Hint: if a pointer may point to two memory regions, these regions *may* be of the same data structure.

Unified Access Hint: if an instruction may access two memory regions, they *may* be of the same data structure.

```

01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (...) {
14.         p = &v;
15.     } else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Our Techniques

Observation 1: while existing techniques mostly focus on memory access patterns (i.e., base addresses and offset values) to identify structures, there are many other program behaviors that can serve as hints to recover data structures.

Allocation Hint: a heap-allocated variable is *likely* to be a structure or an array of structure.

Dataflow Hint: two memory regions connected by a direct dataflow *may* be of the same data structure.

Point-to Hint: if a pointer may point to two memory regions, these regions *may* be of the same data structure.

Unified Access Hint: if an instruction may access two memory regions, they *may* be of the same data structure.

both the heap structure `malloc(...)`,
and the stack structure `v` are accessed by
the same instruction.

```

01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Our Techniques

Observation 1: while existing techniques mostly focus on memory access patterns (i.e., base addresses and offset values) to identify structures, there are many other program behaviors that can serve as hints to recover data structures.

Allocation Hint: a heap-allocated variable is *likely* to be a structure or an array of structure.

Dataflow Hint: two memory regions connected by a direct dataflow *may* be of the same data structure.

Point-to Hint: if a pointer may point to two memory regions, these regions *may* be of the same data structure.

Unified Access Hint: if an instruction may access two memory regions, they *may* be of the same data structure.

```

01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (...) {
14.         p = &v;
15.     } else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Our Techniques

Observation 1: while existing techniques mostly focus on memory access patterns (i.e., base addresses and offset values) to identify structures, there are many other program behaviors that can serve as hints to recover data structures.

Allocation Hint: a heap-allocated variable is *likely* to be a structure or an array of structure.

Dataflow Hint: two memory regions connected by a direct dataflow *may* be of the same data structure.

Point-to Hint: if a pointer may point to two memory regions, these regions *may* be of the same data structure.

Unified Access Hint: if an instruction may access two memory regions, they *may* be of the same data structure.

Other Hints: we also developed other hints (e.g., array hints)

```

01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (...) {
14.         p = &v;
15.     } else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

Our Techniques

Observation 1: while existing techniques mostly focus on memory access patterns (i.e., base addresses and offset values) to identify structures, there are many other program behaviors that can serve as hints to recover data structures.

Allocation Hint: a heap-allocated variable is *likely* to be a structure or an array of structure.

Dataflow Hint: two memory regions connected by a direct dataflow *may* be of the same data structure.

Point-to Hint: if a pointer may point to two memory regions, these regions *may* be of the same data structure.

Unified Access Hint: if an instruction may access two memory regions, they *may* be of the same data structure.

Other Hints: we also developed other hints (e.g., array hints)

How to collect hints?

```

01. typedef struct {
02.     long x;
03.     long y;
04. } elem_t;
05.
06. int main() {
07.     if (!rand(1000)) huft_build(...);
08. }
09.
10. void huft_build(...) {
11.     elem_t *p, v = { .x=0, .y=1 };
12.
13.     if (... )
14.         p = &v;
15.     else {
16.         p = malloc(sizeof(elem_t));
17.         *p = v;
18.     }
19.
20.     output(p->x, p->y);
21. }
```

BDA [OOPSLA'19]: Hint Collection

BDA [OOPSLA'19]: Hint Collection

- BDA is a path sampling driven per-path abstract interpretation technique.
 - For a given sampled path, BDA **ignores the path feasibility** and abstract interprets the program following the given path. Hints can be collected during abstract interpretation.
 - BDA uses **precise symbolic values** as it interprets individual paths separately, which makes it different from other abstract interpretation techniques like VSA that merges values across paths.
- BDA uses a sophisticated path sampling algorithm so that the different paths of a program can be sampled uniformly.
 - Simply tossing a fair coin at each predicate leads to a distribution that is substantially biased towards short paths.
 - Uniform sampling allows exploring a lot more long paths.

BDA [OOPSLA'19]: Hint Collection

- BDA is a path sampling driven per-path abstract interpretation technique.
 - For a given sampled path, BDA **ignores the path feasibility** and abstract interprets the program following the given path. Hints can be collected during abstract interpretation.
 - BDA uses **precise symbolic values** as it interprets individual paths separately, which makes it different from other abstract interpretation techniques like VSA that merges values across paths.
- BDA uses a sophisticated path sampling algorithm so that the different paths of a program can be sampled uniformly.
 - Simply tossing a fair coin at each predicate leads to a distribution that is substantially biased towards short paths.
 - Uniform sampling allows exploring a lot more long paths.

BDA [OOPSLA'19]: Hint Collection

- BDA is a path sampling driven per-path abstract interpretation technique.
 - For a given sampled path, BDA **ignores the path feasibility** and abstract interprets the program following the given path. Hints can be collected during abstract interpretation.
 - BDA uses **precise symbolic values** as it interprets individual paths separately, which makes it different from other abstract interpretation techniques like VSA that merges values across paths.
- BDA uses a sophisticated path sampling algorithm so that the different paths of a program can be sampled uniformly.
 - Simply tossing a fair coin at each predicate leads to a distribution that is substantially biased towards short paths.
 - Uniform sampling allows exploring a lot more long paths.

BDA [OOPSLA'19]: Hint Collection

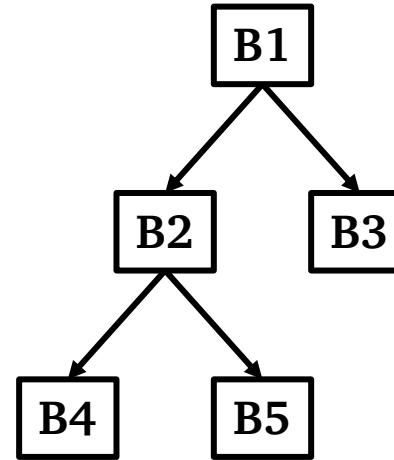
- BDA is a path sampling driven per-path abstract interpretation technique.
 - For a given sampled path, BDA **ignores the path feasibility** and abstract interprets the program following the given path. Hints can be collected during abstract interpretation.
 - BDA uses **precise symbolic values** as it interprets individual paths separately, which makes it different from other abstract interpretation techniques like VSA that merges values across paths.
- BDA uses a sophisticated path sampling algorithm so that the different paths of a program can be sampled uniformly.
 - Simply tossing a fair coin at each predicate leads to a distribution that is substantially biased towards short paths.
 - Uniform sampling allows exploring a lot more long paths.

BDA [OOPSLA'19]: Hint Collection

- BDA is a path sampling driven per-path abstract interpretation technique.
 - For a given sampled path, BDA **ignores the path feasibility** and abstract interprets the program following the given path. Hints can be collected during abstract interpretation.
 - BDA uses **precise symbolic values** as it interprets individual paths separately, which makes it different from other abstract interpretation techniques like VSA that merges values across paths.
- BDA uses a sophisticated path sampling algorithm so that the different paths of a program can be sampled uniformly.
 - Simply tossing a fair coin at each predicate leads to a distribution that is substantially biased towards short paths.
 - Uniform sampling allows exploring a lot more long paths.

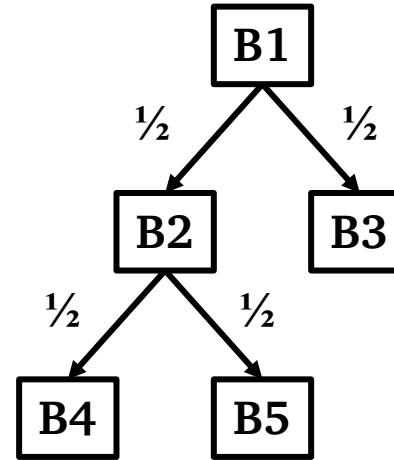
BDA [OOPSLA'19]: Hint Collection

- BDA is a path sampling driven per-path abstract interpretation technique.
 - For a given sampled path, BDA **ignores the path feasibility** and abstract interprets the program following the given path. Hints can be collected during abstract interpretation.
 - BDA uses **precise symbolic values** as it interprets individual paths separately, which makes it different from other abstract interpretation techniques like VSA that merges values across paths.
- BDA uses a sophisticated path sampling algorithm so that the different paths of a program can be sampled uniformly.
 - Simply tossing a fair coin at each predicate leads to a distribution that is substantially biased towards short paths.
 - Uniform sampling allows exploring a lot more long paths.



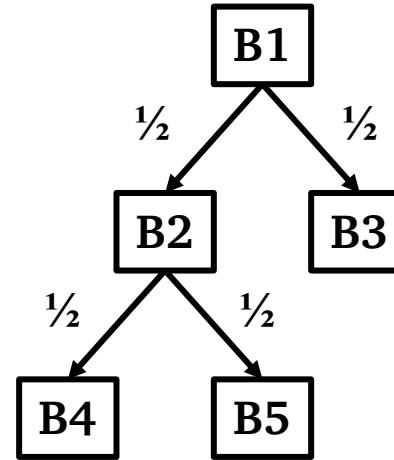
BDA [OOPSLA'19]: Hint Collection

- BDA is a path sampling driven per-path abstract interpretation technique.
 - For a given sampled path, BDA **ignores the path feasibility** and abstract interprets the program following the given path. Hints can be collected during abstract interpretation.
 - BDA uses **precise symbolic values** as it interprets individual paths separately, which makes it different from other abstract interpretation techniques like VSA that merges values across paths.
- BDA uses a sophisticated path sampling algorithm so that the different paths of a program can be sampled uniformly.
 - Simply tossing a fair coin at each predicate leads to a distribution that is substantially biased towards short paths.
 - Uniform sampling allows exploring a lot more long paths.



BDA [OOPSLA'19]: Hint Collection

- BDA is a path sampling driven per-path abstract interpretation technique.
 - For a given sampled path, BDA **ignores the path feasibility** and abstract interprets the program following the given path. Hints can be collected during abstract interpretation.
 - BDA uses **precise symbolic values** as it interprets individual paths separately, which makes it different from other abstract interpretation techniques like VSA that merges values across paths.
- BDA uses a sophisticated path sampling algorithm so that the different paths of a program can be sampled uniformly.
 - Simply tossing a fair coin at each predicate leads to a distribution that is substantially biased towards short paths.
 - Uniform sampling allows exploring a lot more long paths.



Naïve Sampling:

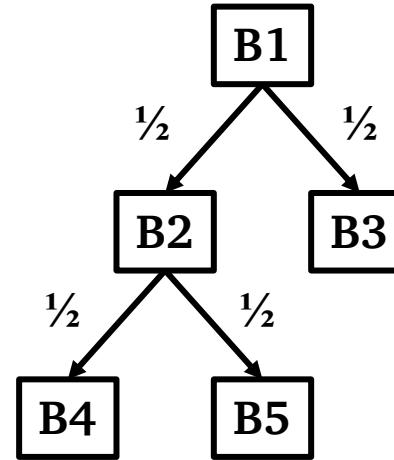
$$\mathbf{B1 - B3: } \frac{1}{2}$$

$$\mathbf{B1 - B2 - B4: } \frac{1}{4}$$

$$\mathbf{B1 - B2 - B5: } \frac{1}{4}$$

BDA [OOPSLA'19]: Hint Collection

- BDA is a path sampling driven per-path abstract interpretation technique.
 - For a given sampled path, BDA **ignores the path feasibility** and abstract interprets the program following the given path. Hints can be collected during abstract interpretation.
 - BDA uses **precise symbolic values** as it interprets individual paths separately, which makes it different from other abstract interpretation techniques like VSA that merges values across paths.
- BDA uses a sophisticated path sampling algorithm so that the different paths of a program can be sampled uniformly.
 - Simply tossing a fair coin at each predicate leads to a distribution that is substantially biased towards short paths.
 - Uniform sampling allows exploring a lot more long paths.



Naïve Sampling:

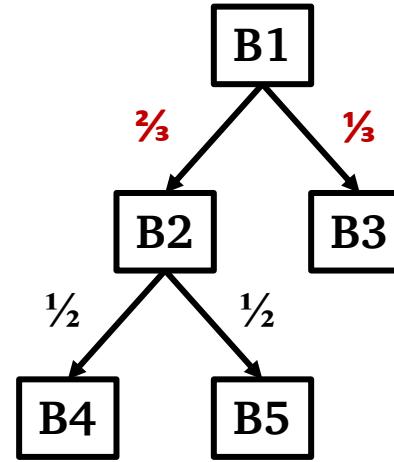
$$\mathbf{B1 - B3: } \frac{1}{2}$$

$$\mathbf{B1 - B2 - B4: } \frac{1}{4}$$

$$\mathbf{B1 - B2 - B5: } \frac{1}{4}$$

BDA [OOPSLA'19]: Hint Collection

- BDA is a path sampling driven per-path abstract interpretation technique.
 - For a given sampled path, BDA **ignores the path feasibility** and abstract interprets the program following the given path. Hints can be collected during abstract interpretation.
 - BDA uses **precise symbolic values** as it interprets individual paths separately, which makes it different from other abstract interpretation techniques like VSA that merges values across paths.
- BDA uses a sophisticated path sampling algorithm so that the different paths of a program can be sampled uniformly.
 - Simply tossing a fair coin at each predicate leads to a distribution that is substantially biased towards short paths.
 - Uniform sampling allows exploring a lot more long paths.



Naïve Sampling:

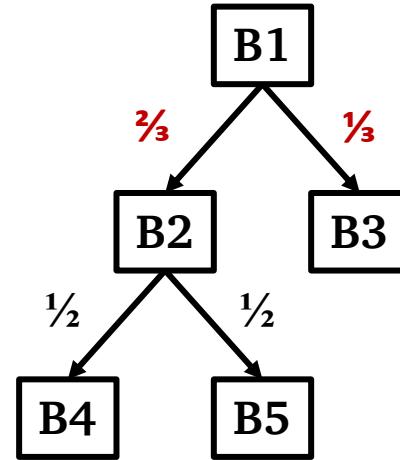
$$\mathbf{B1 - B3: } \frac{1}{2}$$

$$\mathbf{B1 - B2 - B4: } \frac{1}{4}$$

$$\mathbf{B1 - B2 - B5: } \frac{1}{4}$$

BDA [OOPSLA'19]: Hint Collection

- BDA is a path sampling driven per-path abstract interpretation technique.
 - For a given sampled path, BDA **ignores the path feasibility** and abstract interprets the program following the given path. Hints can be collected during abstract interpretation.
 - BDA uses **precise symbolic values** as it interprets individual paths separately, which makes it different from other abstract interpretation techniques like VSA that merges values across paths.
- BDA uses a sophisticated path sampling algorithm so that the different paths of a program can be sampled uniformly.
 - Simply tossing a fair coin at each predicate leads to a distribution that is substantially biased towards short paths.
 - Uniform sampling allows exploring a lot more long paths.



Naïve Sampling:

$$\mathbf{B1 - B3: } \frac{1}{2}$$

$$\mathbf{B1 - B2 - B4: } \frac{1}{4}$$

$$\mathbf{B1 - B2 - B5: } \frac{1}{4}$$

BDA Sampling:

$$\mathbf{B1 - B3: } \frac{1}{3}$$

$$\mathbf{B1 - B2 - B4: } \frac{1}{3}$$

$$\mathbf{B1 - B2 - B5: } \frac{1}{3}$$

Our Techniques

Observation 1: while existing techniques mostly focus on memory access patterns (i.e., base addresses and offset values) to identify structures, there are many other program behaviors that can serve as hints to recover data structures.

Allocation Hint: a heap-allocated variable is *likely* to be a structure or an array of structure.

Dataflow Hint: two memory regions connected by a direct dataflow *may* be of the same data structure.

Point-to Hint: if a pointer may point to two memory regions, these regions *may* be of the same data structure.

Unified Access Hint: if an instruction may access two memory regions, they *may* be of the same data structure.

Other Hints: we also developed other hints (e.g., array hints)

Our Techniques

Observation 1: while existing techniques mostly focus on memory access patterns (i.e., base addresses and offset values) to identify structures, there are many other program behaviors that can serve as hints to recover data structures.

Allocation Hint: a heap-allocated variable is *likely* to be a structure or an array of structure.

Dataflow Hint: two memory regions connected by a direct dataflow *may* be of the same data structure.

Point-to Hint: if a pointer may point to two memory regions, these regions *may* be of the same data structure.

Unified Access Hint: if an instruction may access two memory regions, they *may* be of the same data structure.

Other Hints: we also developed other hints (e.g., array hints)

Uncertainty

Our Techniques

Observation 2: the various kinds of hints in variable/structure recovery can be integrated in a more sophisticated manner using **probabilistic inference**.

Our Techniques

Observation 2: the various kinds of hints in variable/structure recovery can be integrated in a more sophisticated manner using **probabilistic inference**.

AllocationHint

PointToHint

DataFlowHint

UnifiedAccessHint

Our Techniques

Observation 2: the various kinds of hints in variable/structure recovery can be integrated in a more sophisticated manner using **probabilistic inference**.

AllocationHint 0.9

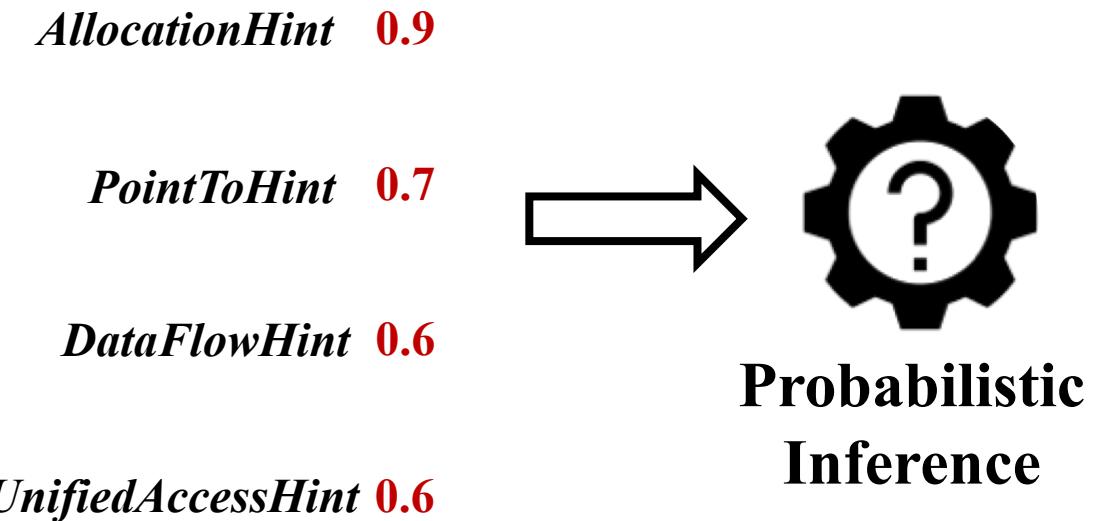
PointToHint 0.7

DataFlowHint 0.6

UnifiedAccessHint 0.6

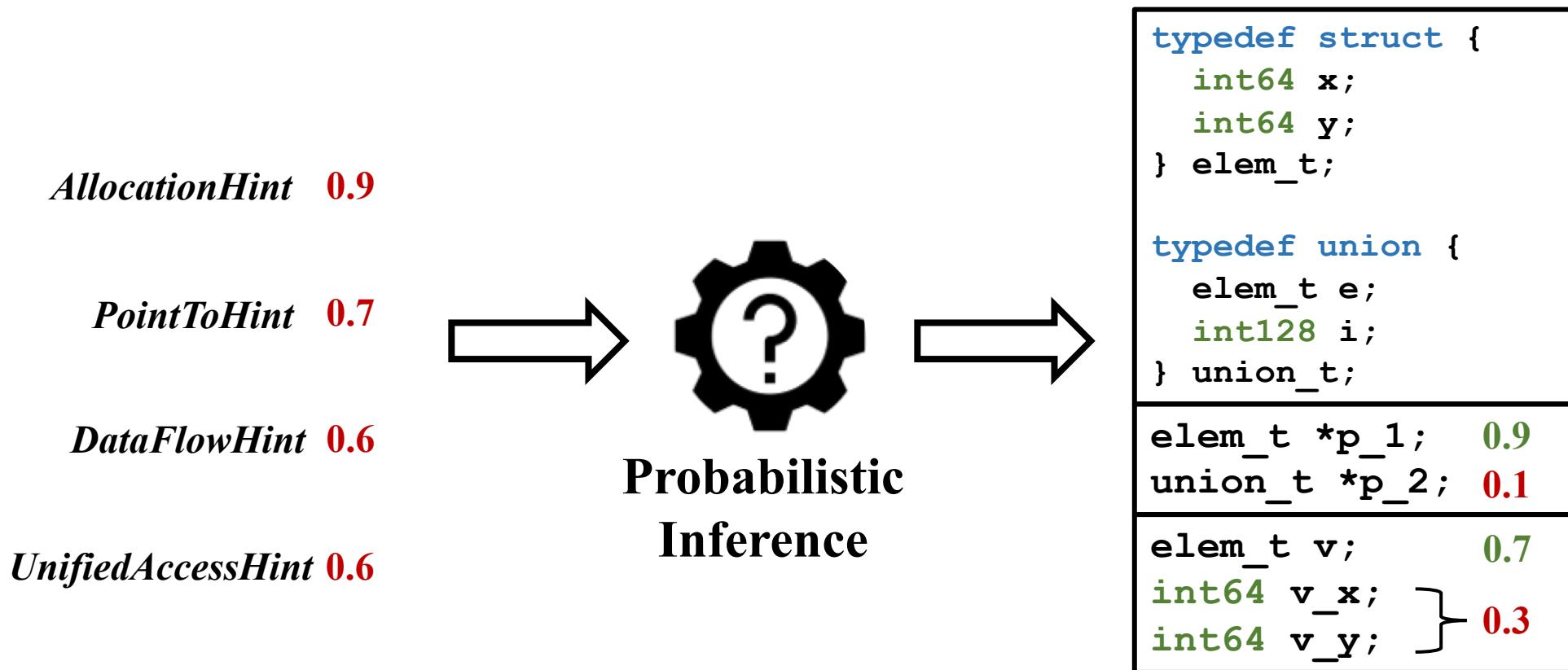
Our Techniques

Observation 2: the various kinds of hints in variable/structure recovery can be integrated in a more sophisticated manner using **probabilistic inference**.



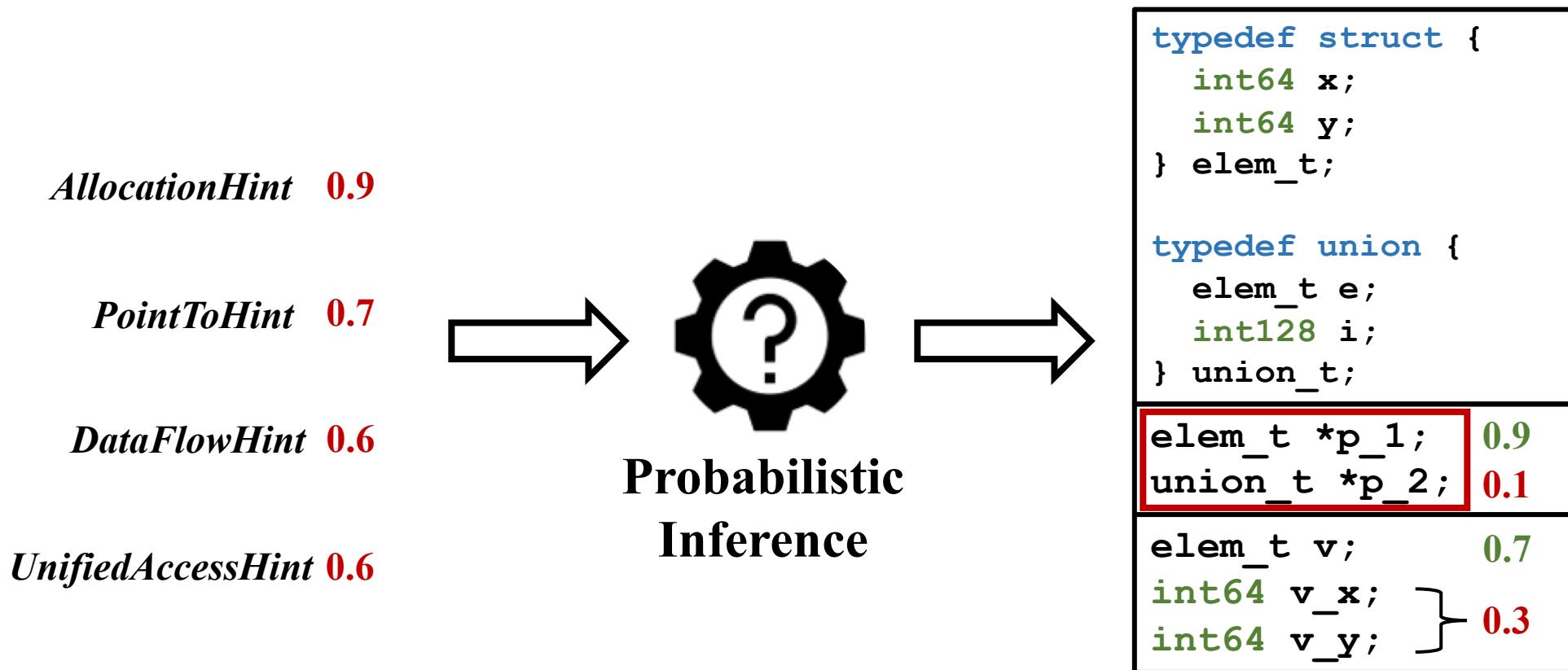
Our Techniques

Observation 2: the various kinds of hints in variable/structure recovery can be integrated in a more sophisticated manner using **probabilistic inference**.



Our Techniques

Observation 2: the various kinds of hints in variable/structure recovery can be integrated in a more sophisticated manner using **probabilistic inference**.



Our Techniques

Observation 2: the various kinds of hints in variable/structure recovery can be integrated in a more sophisticated manner using **probabilistic inference**.

Our Techniques

Observation 2: the various kinds of hints in variable/structure recovery can be integrated in a more sophisticated manner using **probabilistic inference**.

Random Variables: Random variables are introduced to describe variable properties, type properties, and structural properties.

- They are comprehensive and can be used to express the needed symbol information

Our Techniques

Observation 2: the various kinds of hints in variable/structure recovery can be integrated in a more sophisticated manner using **probabilistic inference**.

Random Variables: Random variables are introduced to describe variable properties, type properties, and structural properties.

- They are comprehensive and can be used to express the needed symbol information

Our Techniques

Observation 2: the various kinds of hints in variable/structure recovery can be integrated in a more sophisticated manner using **probabilistic inference**.

Random Variables: Random variables are introduced to describe variable properties, type properties, and structural properties.

- They are comprehensive and can be used to express the needed symbol information

PrimitiveVar(a, s): a memory region at address **a** with size **s** is a primitive variable.

Our Techniques

Observation 2: the various kinds of hints in variable/structure recovery can be integrated in a more sophisticated manner using **probabilistic inference**.

Random Variables: Random variables are introduced to describe variable properties, type properties, and structural properties.

- They are comprehensive and can be used to express the needed symbol information

PrimitiveVar(a, s): a memory region at address **a** with size **s** is a primitive variable.

FieldOf(a1, s, a2): a memory region at address **a1** with size **s** is a field of a structure at address **a2**.

Our Techniques

Observation 2: the various kinds of hints in variable/structure recovery can be integrated in a more sophisticated manner using **probabilistic inference**.

Random Variables: Random variables are introduced to describe variable properties, type properties, and structural properties.

- They are comprehensive and can be used to express the needed symbol information

PrimitiveVar(a, s): a memory region at address **a** with size **s** is a primitive variable.

FieldOf(a1, s, a2): a memory region at address **a1** with size **s** is a field of a structure at address **a2**.

HomoRegion(a1, a2, s): the memory regions starting at address **a1** and **a2** (with size **s**) have homogeneous structures.

Our Techniques

Observation 2: the various kinds of hints in variable/structure recovery can be integrated in a more sophisticated manner using **probabilistic inference**.

Random Variables: Random variables are introduced to describe variable properties, type properties, and structural properties.

PrimitiveVar(a, s)

FieldOf(a1, s, a2)

HomoRegion(a1, a2, s)

Our Techniques

Observation 2: the various kinds of hints in variable/structure recovery can be integrated in a more sophisticated manner using **probabilistic inference**.

Random Variables: Random variables are introduced to describe variable properties, type properties, and structural properties.

PrimitiveVar(a, s)

FieldOf(a1, s, a2)

HomoRegion(a1, a2, s)

Probabilistic Constraints: Probabilistic constraints essentially denote probability functions over the random variables involved.

- Correlate these random variables
- Calculate the posterior marginal probabilities for random variables
- The inference denoted by a constraint is probabilistic

Our Techniques

Observation 2: the various kinds of hints in variable/structure recovery can be integrated in a more sophisticated manner using **probabilistic inference**.

Random Variables: Random variables are introduced to describe variable properties, type properties, and structural properties.

PrimitiveVar(a, s)

FieldOf(a1, s, a2)

HomoRegion(a1, a2, s)

Probabilistic Constraints: Probabilistic constraints essentially denote probability functions over the random variables involved.

- Correlate these random variables
- Calculate the posterior marginal probabilities for random variables
- The inference denoted by a constraint is probabilistic

Our Techniques

Observation 2: the various kinds of hints in variable/structure recovery can be integrated in a more sophisticated manner using **probabilistic inference**.

Random Variables: Random variables are introduced to describe variable properties, type properties, and structural properties.

PrimitiveVar(a, s)

FieldOf(a1, s, a2)

HomoRegion(a1, a2, s)

Probabilistic Constraints: Probabilistic constraints essentially denote probability functions over the random variables involved.

- Correlate these random variables
- Calculate the posterior marginal probabilities for random variables
- The inference denoted by a constraint is probabilistic

Our Techniques

Observation 2: the various kinds of hints in variable/structure recovery can be integrated in a more sophisticated manner using **probabilistic inference**.

Random Variables: Random variables are introduced to describe variable properties, type properties, and structural properties.

PrimitiveVar(a, s)

FieldOf(a1, s, a2)

HomoRegion(a1, a2, s)

Probabilistic Constraints: Probabilistic constraints essentially denote probability functions over the random variables involved.

- Correlate these random variables
- Calculate the posterior marginal probabilities for random variables
- The inference denoted by a constraint is probabilistic

Our Techniques

Observation 2: the various kinds of hints in variable/structure recovery can be integrated in a more sophisticated manner using **probabilistic inference**.

Random Variables: Random variables are introduced to describe variable properties, type properties, and structural properties.

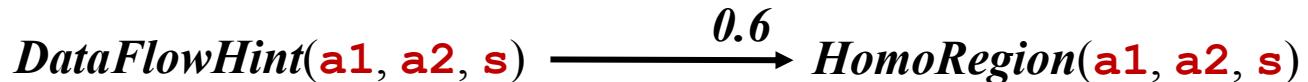
PrimitiveVar(a, s)

FieldOf(a1, s, a2)

HomoRegion(a1, a2, s)

Probabilistic Constraints: Probabilistic constraints essentially denote probability functions over the random variables involved.

- Correlate these random variables
- Calculate the posterior marginal probabilities for random variables
- The inference denoted by a constraint is probabilistic



Our Techniques

Observation 2: the various kinds of hints in variable/structure recovery can be integrated in a more sophisticated manner using **probabilistic inference**.

Random Variables: Random variables are introduced to describe variable properties, type properties, and structural properties.

$\text{PrimitiveVar}(\mathbf{a}, \mathbf{s})$

$\text{FieldOf}(\mathbf{a1}, \mathbf{s}, \mathbf{a2})$

$\text{HomoRegion}(\mathbf{a1}, \mathbf{a2}, \mathbf{s})$

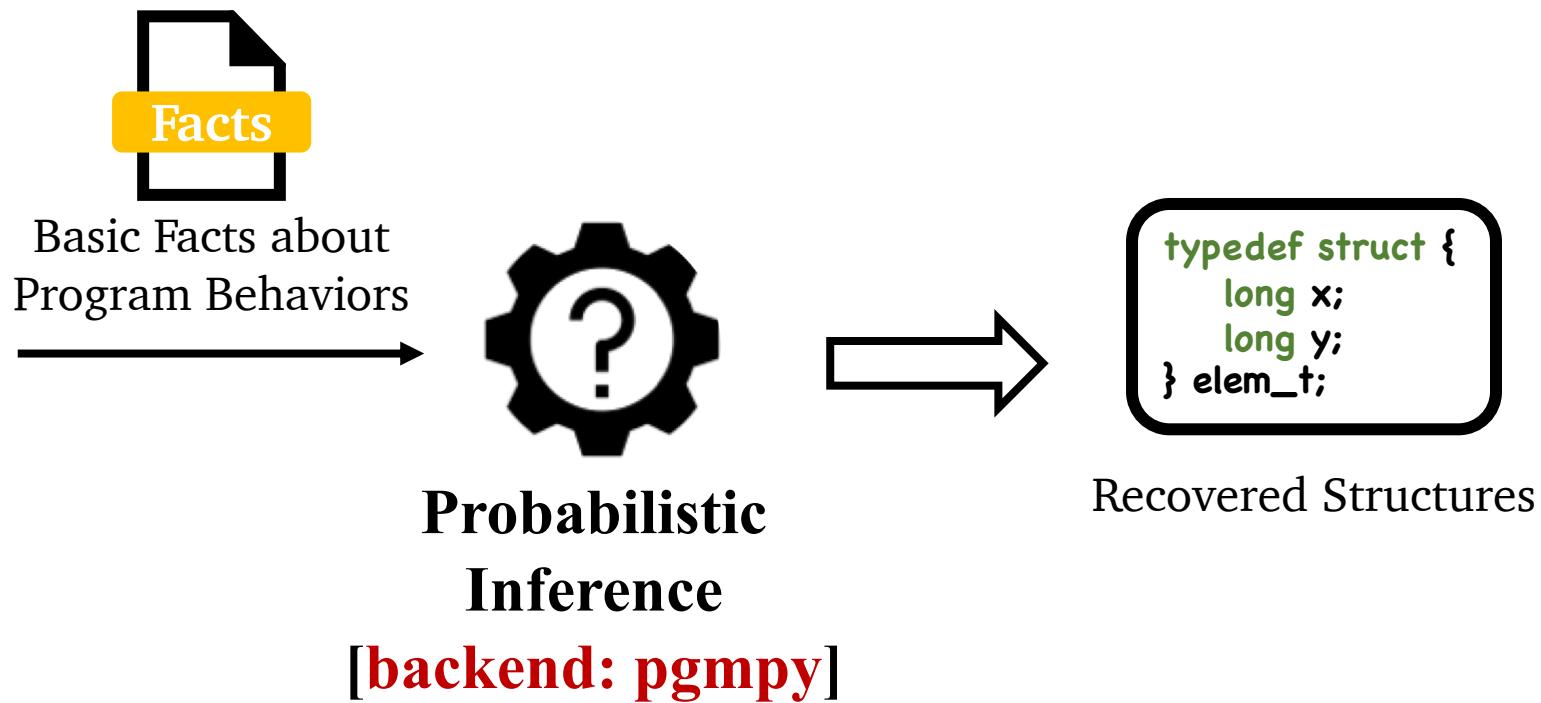
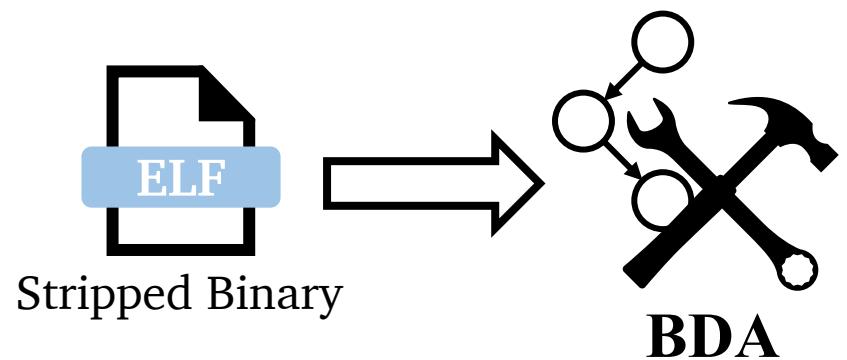
Probabilistic Constraints: Probabilistic constraints essentially denote probability functions over the random variables involved.

- Correlate these random variables
- Calculate the posterior marginal probabilities for random variables
- The inference denoted by a constraint is probabilistic

$$\text{DataFlowHint}(\mathbf{a1}, \mathbf{a2}, \mathbf{s}) \xrightarrow{0.6} \text{HomoRegion}(\mathbf{a1}, \mathbf{a2}, \mathbf{s})$$

$$\text{FieldOf}(\mathbf{a1}, \mathbf{s}, \mathbf{a2}) \wedge \text{HomoRegion}(\mathbf{a2}, \mathbf{a3}, \underline{\mathbf{s}}) \xrightarrow{0.9} \text{FieldOf}(\mathbf{a3} + (\mathbf{a1} - \mathbf{a2}), \mathbf{s}, \mathbf{a3})$$

Osprey 



Evaluation

Evaluation

Benchmark:

- CoreUtils (O0 and O3)
- Benchmark used by Howard Project (O0 and O3)

Ground Truth:

- debugging information

Evaluation

Benchmark:

- CoreUtils (O0 and O3)
- Benchmark used by Howard Project (O0 and O3)

Ground Truth:

- debugging information

Evaluation

Benchmark:

- CoreUtils (O0 and O3)
- Benchmark used by Howard Project (O0 and O3)

Ground Truth:

- debugging information

Evaluation

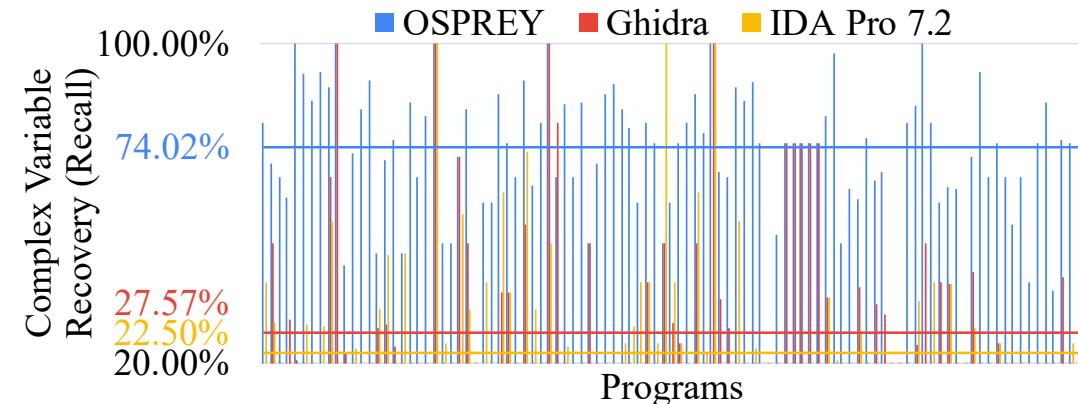
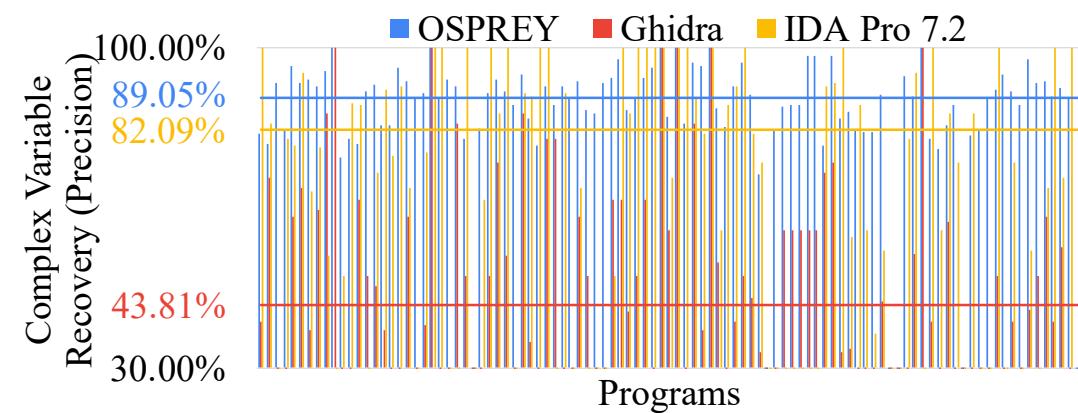
Benchmark:

- CoreUtils (O0 and O3)
- Benchmark used by Howard Project (O0 and O3)

Ground Truth:

- debugging information

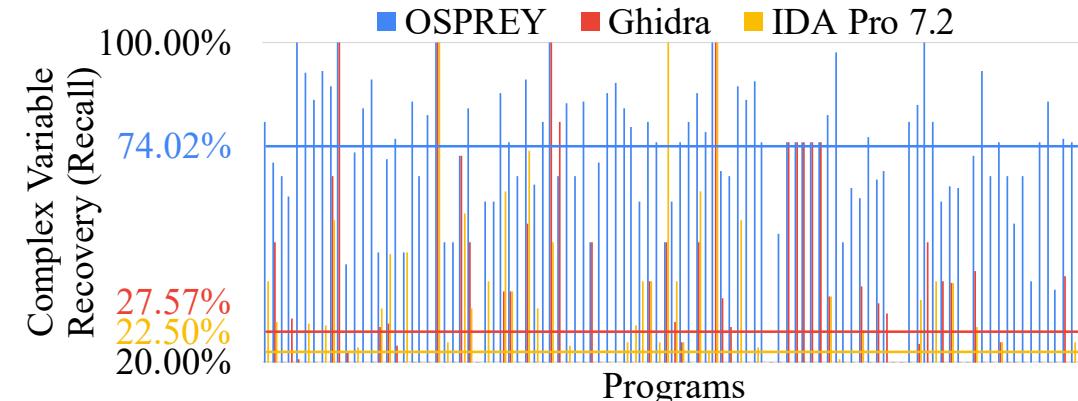
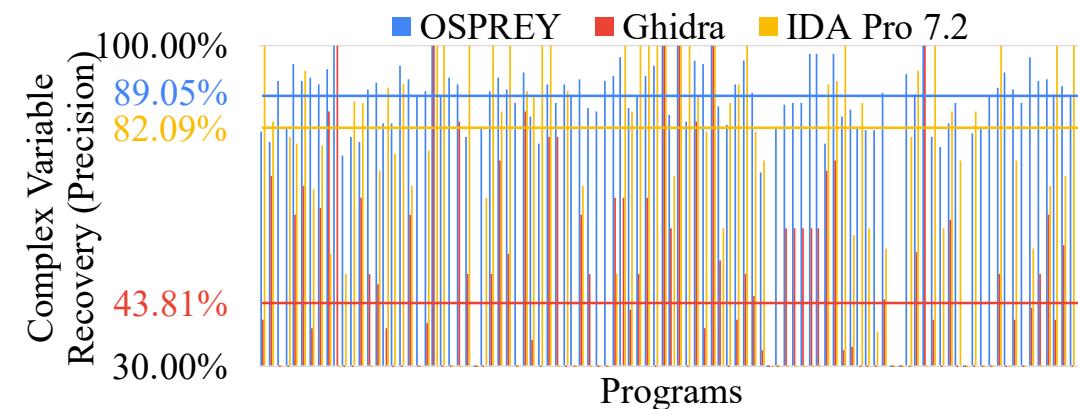
Evaluation: Complex Variable Recovery on CoreUtils



Evaluation: Complex Variable Recovery on CoreUtils

Similar to the standard in the literature, we inspect individual variables.

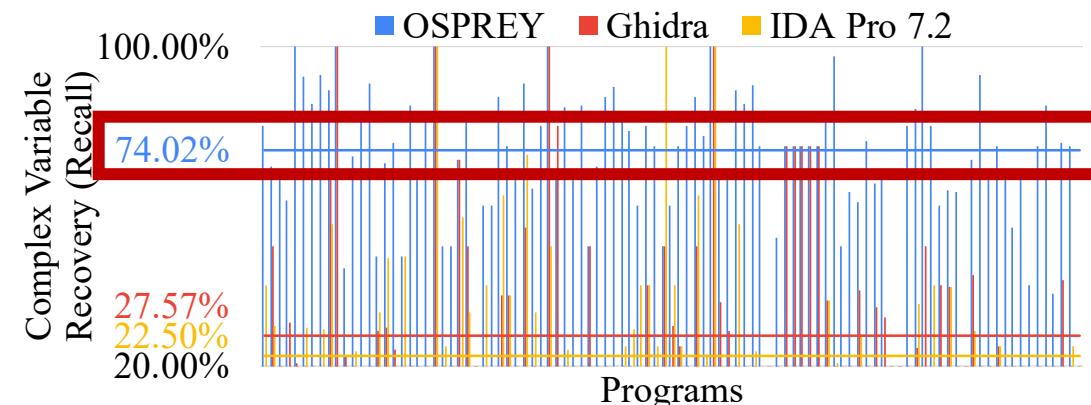
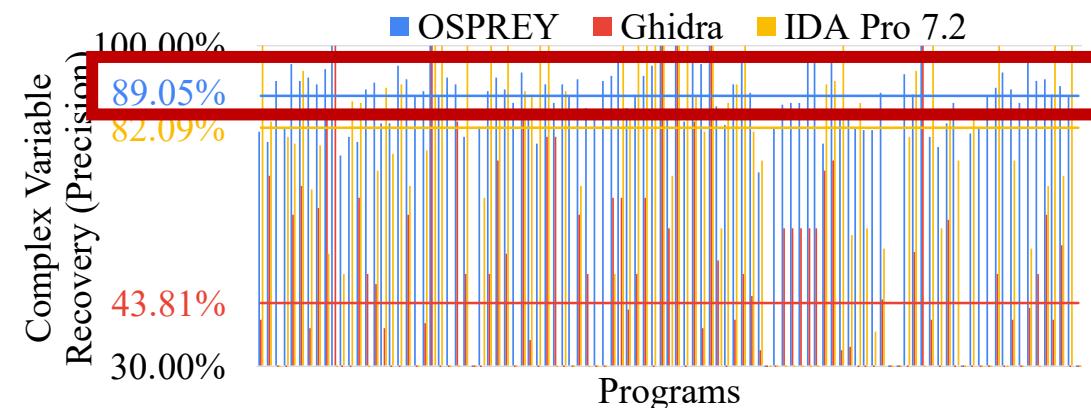
- Specially, if it is a pointer type, we inspect the structure that is being pointed to.
- Complex variables include structures, unions, arrays and pointers to structures and unions.



Evaluation: Complex Variable Recovery on CoreUtils

Similar to the standard in the literature, we inspect individual variables.

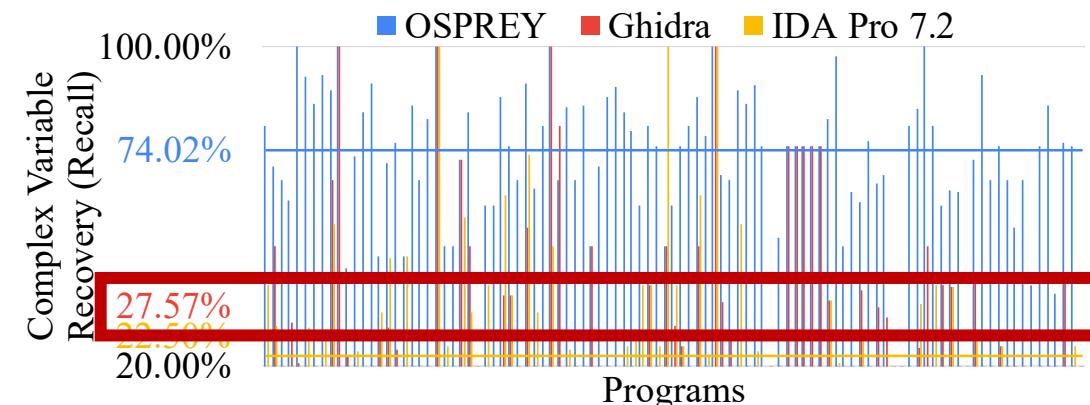
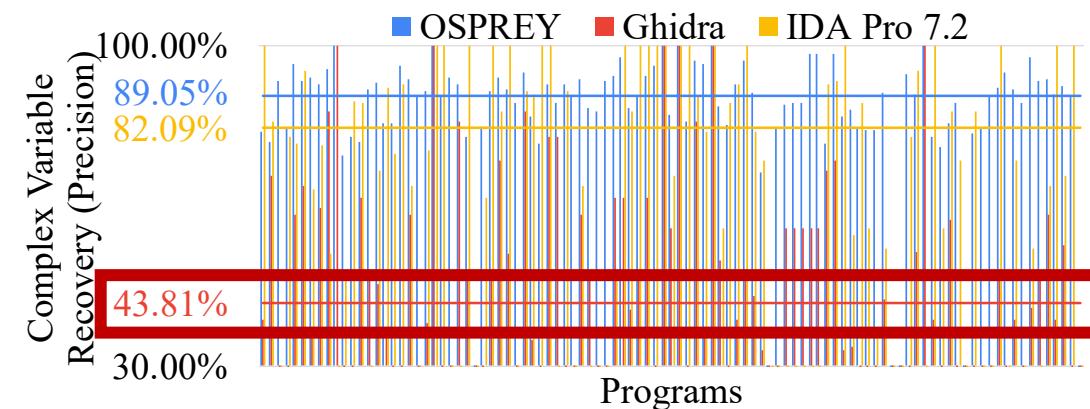
- Specially, if it is a pointer type, we inspect the structure that is being pointed to.
- Complex variables include structures, unions, arrays and pointers to structures and unions.



Evaluation: Complex Variable Recovery on CoreUtils

Similar to the standard in the literature, we inspect individual variables.

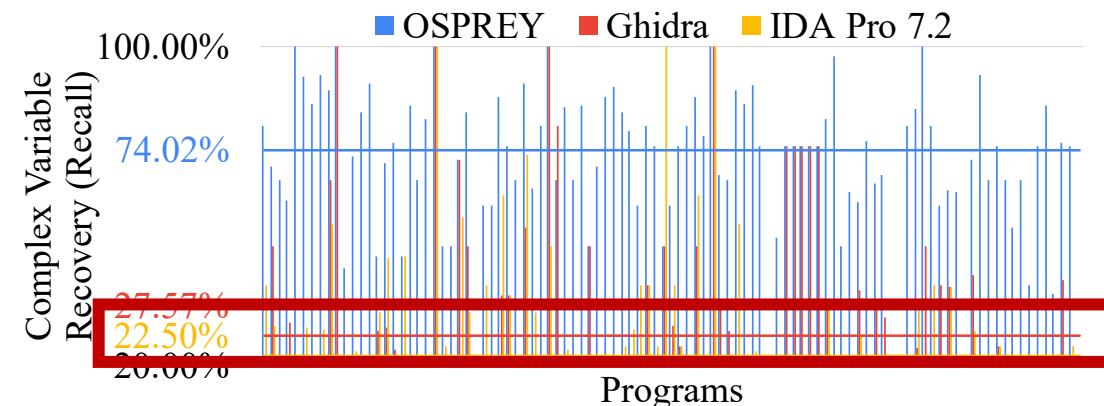
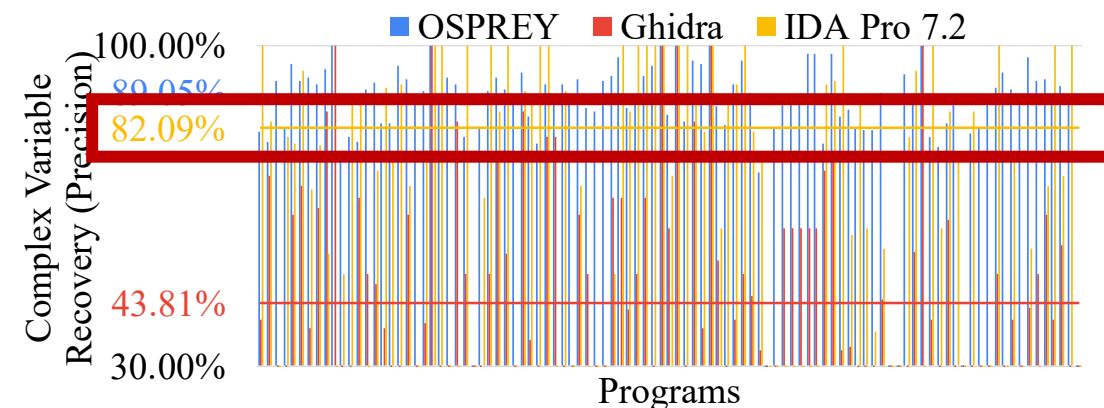
- Specially, if it is a pointer type, we inspect the structure that is being pointed to.
- Complex variables include structures, unions, arrays and pointers to structures and unions.



Evaluation: Complex Variable Recovery on CoreUtils

Similar to the standard in the literature, we inspect individual variables.

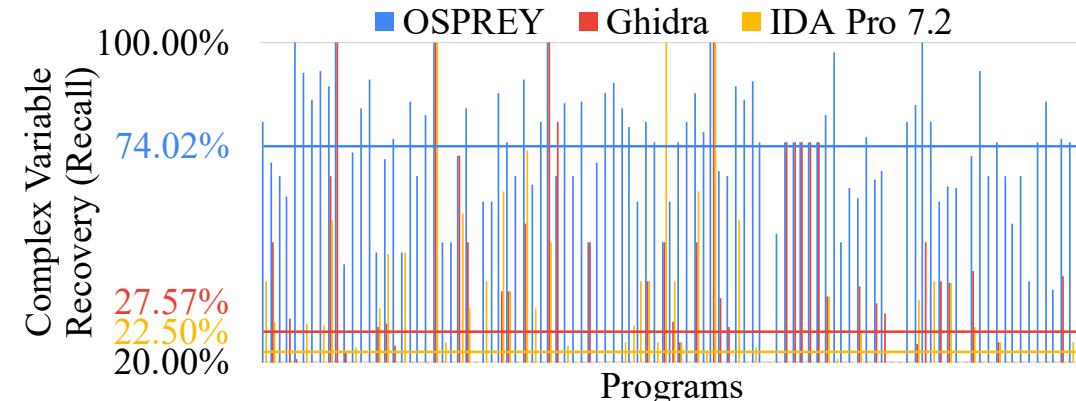
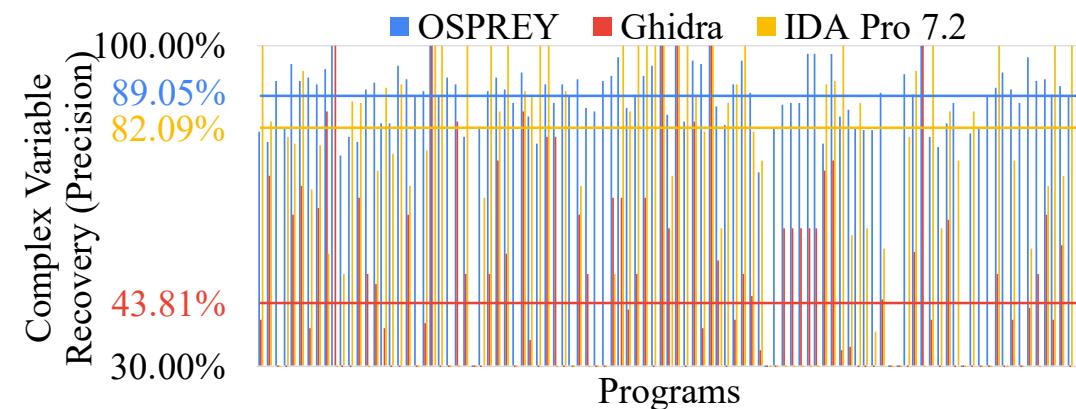
- Specially, if it is a pointer type, we inspect the structure that is being pointed to.
- Complex variables include structures, unions, arrays and pointers to structures and unions.



Evaluation: Complex Variable Recovery on CoreUtils

Similar to the standard in the literature, we inspect individual variables.

- Specially, if it is a pointer type, we inspect the structure that is being pointed to.
- Complex variables include structures, unions, arrays and pointers to structures and unions.

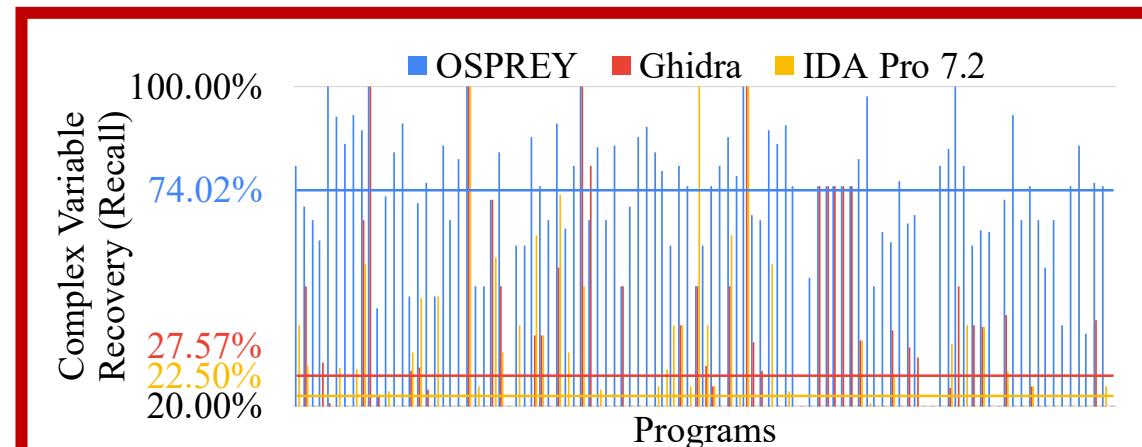
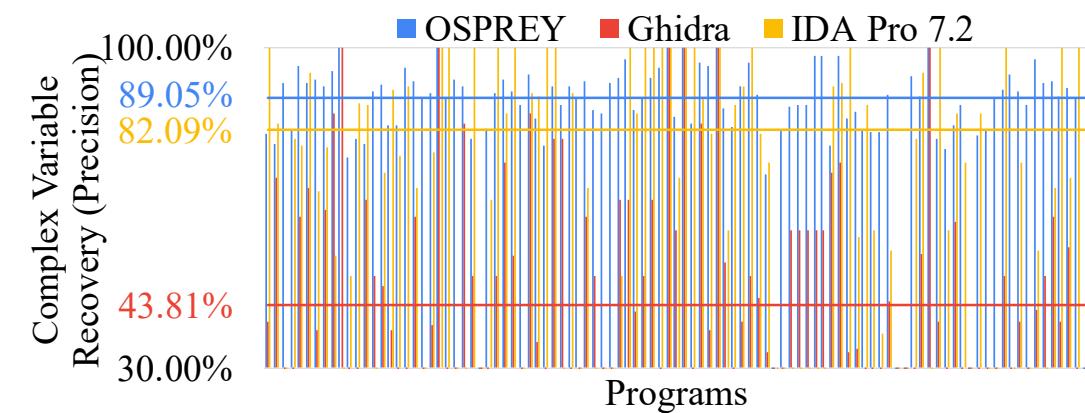


Evaluation: Complex Variable Recovery on CoreUtils

Similar to the standard in the literature, we inspect individual variables.

- Specially, if it is a pointer type, we inspect the structure that is being pointed to.
- Complex variables include structures, unions, arrays and pointers to structures and unions.

- The recall of OSPREY is around **74%**, more than **2** times higher than Ghidra and IDA Pro.
- The precision of OSPREY also outperforms Ghidra and IDA Pro.
 - We inspect why such improvement, and the main reason is that IDA and Ghidra do not perform complex analysis on recover structure while most complex variables are structure and pointers to structures.

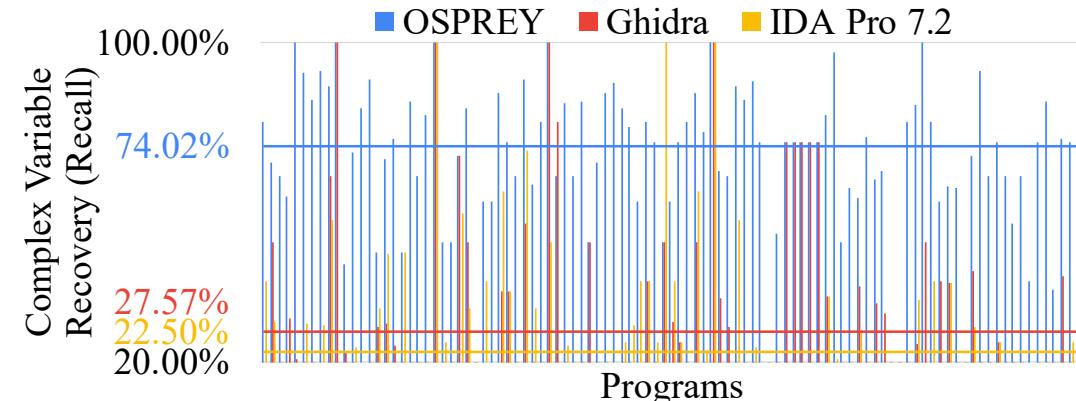
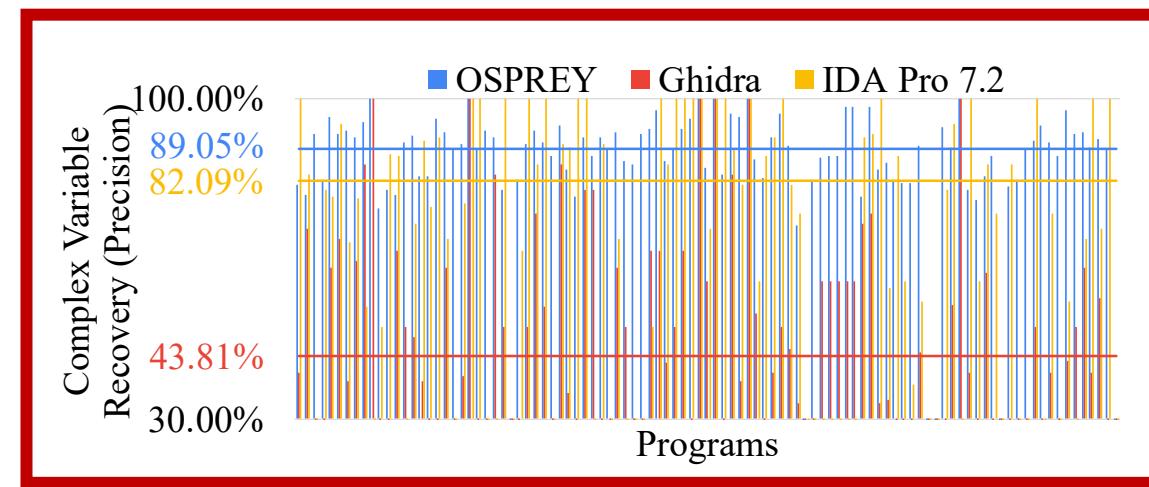


Evaluation: Complex Variable Recovery on CoreUtils

Similar to the standard in the literature, we inspect individual variables.

- Specially, if it is a pointer type, we inspect the structure that is being pointed to.
- Complex variables include structures, unions, arrays and pointers to structures and unions.

- The recall of OSPREY is around **74%**, more than **2** times higher than Ghidra and IDA Pro.
- The precision of OSPREY also outperforms Ghidra and IDA Pro.
 - We inspect why such improvement, and the main reason is that IDA and Ghidra do not perform complex analysis on recover structure while most complex variables are structure and pointers to structures.

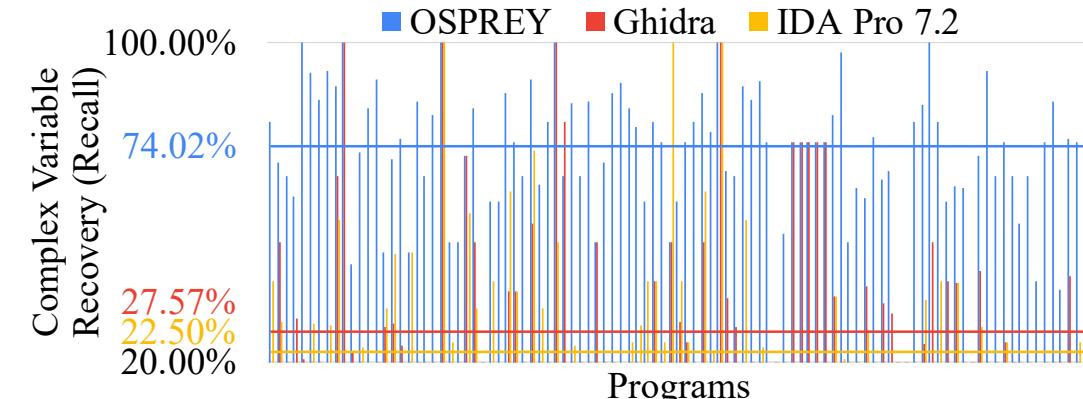
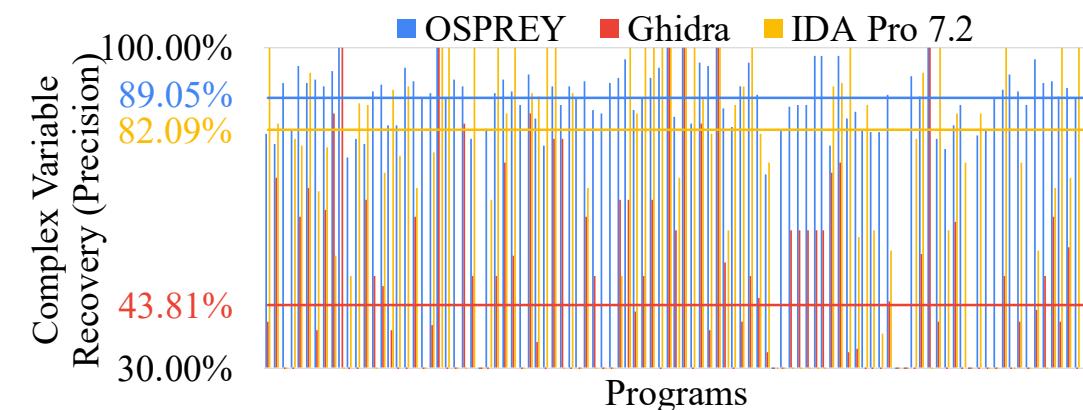


Evaluation: Complex Variable Recovery on CoreUtils

Similar to the standard in the literature, we inspect individual variables.

- Specially, if it is a pointer type, we inspect the structure that is being pointed to.
- Complex variables include structures, unions, arrays and pointers to structures and unions.

- The recall of OSPREY is around **74%**, more than **2** times higher than Ghidra and IDA Pro.
- The precision of OSPREY also outperforms Ghidra and IDA Pro.
 - We inspect why such improvement, and the main reason is that IDA and Ghidra do not perform complex analysis on recover structure while most complex variables are structure and pointers to structures.



Evaluation: Argument Decomposition Results

```

1 int network_rxxx(server *srv)
2 {
3     server *v1; // rbx
4     int result; // rax
5     size_t v3; // rbp
6     server_socket *v4; // r12
7     fdnode *v5; // rax
8     fdevents *v6; // rdi
9
10    v1 = srv;
11    result = fdevent_sxxx(srv->ev);
12    if ( result != -1 )
13    {
14        v3 = 0LL;
15        if ( !srv->sockets_disabled )
16        {
17            while ( v1->srv_sockets.
18                  used > v3 )
19            {
20                v4 = v1->srv_sockets.
21                      ptr[v3++];
22                v5 = fdevent_gxxx(
23                    v1->ev,
24                    v4->fd,
25                    network_sxxx, v4
26                );
27                v6 = v1->ev;
28                v4->fdn = v5;
29                fdevent_fxxx(v6, v5, 1);
30            }
31        }
32        result = 0LL;
33    }
34    return result;
35 }
```

IDA Pro w/ debug information

```

1 __int32 __fastcall sub_D840(__int64 a1)
2 {
3     _QWORD *v1; // rbx
4     __int32 result; // rax
5     unsigned __int64 v3; // rbp
6     __int64 v4; // r12
7     __int64 v5; // rax
8     __int64 v6; // rdi
9
10    v1 = (_QWORD *)a1;
11    result = sub_12B7A(*(_QWORD *)(a1 + 24));
12    if ( (_DWORD)result != -1 )
13    {
14        v3 = 0LL;
15        if ( !*(__DWORD *)(a1 + 100) )
16        {
17            while ( v1[2] > v3 )
18            {
19                v4 = *(_QWORD *)(*v1 + 8 * v3++);
20
21                v5 = sub_21860(
22                    v1[3],
23                    *(unsigned int *) (v4 + 112),
24                    sub_18F30, v4
25                );
26                v6 = v1[3];
27                *(_QWORD *) (v4 + 120) = v5;
28                sub_219C0(v6, v5, 1);
29            }
30        }
31        result = 0;
32    }
33    return result;
34 }
```

Vanilla IDA Pro

```

1 __int32 __fastcall sub_D840(struct_C264 *a1)
2 {
3     struct_C264 *v1; // rbx
4     __int32 result; // rax
5     unsigned __int64 v3; // rbp
6     struct_CF4A *v4; // r12
7     struct_12A42 *v5; // rax
8     struct_12A0E *v6; // rdi
9
10    v1 = a1;
11    result = sub_12B7A(a1->ptr_field_28);
12    if ( result != -1 )
13    {
14        v3 = 0LL;
15        if ( !a1->dat_field_74 )
16        {
17            while ( v1->dat_field_10 > v3 )
18            {
19                v4 = v1->ptr_ptr_field_0[v3++];
20
21                v5 = sub_21860(
22                    v1->ptr_field_28,
23                    v4->dat_field_10,
24                    sub_18F30, v4
25                );
26                v6 = v1->ptr_field_28;
27                v4->ptr_field_18 = v5;
28                sub_219C0(v6, v5, 1);
29            }
30        }
31        result = 0;
32    }
33    return result;
34 }
```

IDA Pro w/ our tool

Evaluation: Argument Decomposition Results

```

1 int network_rxxx(server *srv)
2 {
3     server *v1; // rbx
4     int result; // rax
5     size_t v3; // rbp
6     server_socket *v4; // r12
7     fnode *v5; // rax
8     fdevents *v6; // rdi
9
10    v1 = srv;
11    result = fdevent_sxxx(srv->ev);
12    if ( result != -1 )
13    {
14        v3 = 0LL;
15        if ( !srv->sockets_disabled )
16        {
17            while ( v1->srv_sockets.
18                  used > v3 )
19            {
20                v4 = v1->srv_sockets.
21                      ptr[v3++];
22                v5 = fdevent_gxxx(
23                    v1->ev,
24                    v4->fd,
25                    network_sxxx, v4
26                );
27                v6 = v1->ev;
28                v4->fdn = v5;
29                fdevent_fxxx(v6, v5, 1);
30            }
31        }
32        result = 0LL;
33    }
34    return result;
35 }
```

IDA Pro w/ debug information

```

1 __int32 __fastcall sub_D840(__int64 a1)
2 {
3     QWORD *v1; // rbx
4     __int32 result; // rax
5     unsigned __int64 v3; // rbp
6     __int64 v4; // r12
7     __int64 v5; // rax
8     __int64 v6; // rdi
9
10    v1 = (_QWORD *)a1;
11    result = sub_12B7A(*(_QWORD *)(a1 + 24));
12    if ( (_DWORD)result != -1 )
13    {
14        v3 = 0LL;
15        if ( !*(DWORD *)(a1 + 100) )
16        {
17            while ( v1[2] > v3 )
18            {
19                v4 = *(_QWORD *)(*v1 + 8 * v3++);
20                v5 = sub_21860(
21                    v1[3],
22                    *(unsigned int *)(v4 + 112),
23                    sub_18F30, v4
24                );
25                v6 = v1[3];
26                (*(_QWORD *) (v4 + 120)) = v5;
27                sub_219C0(v6, v5, 1);
28            }
29            result = 0;
30        }
31    }
32    return result;
33 }
```

Vanilla IDA Pro

```

1 __int32 __fastcall sub_D840(struct_C264 *a1)
2 {
3     struct_C264 *v1; // rbx
4     __int32 result; // rax
5     unsigned __int64 v3; // rbp
6     struct_CF4A *v4; // r12
7     struct_12A42 *v5; // rax
8     struct_12A0E *v6; // rdi
9
10    v1 = a1;
11    result = sub_12B7A(a1->ptr_field_28);
12    if ( result != -1 )
13    {
14        v3 = 0LL;
15        if ( !a1->dat_field_74 )
16        {
17            while ( v1->dat_field_10 > v3 )
18            {
19                v4 = v1->ptr_ptr_field_0[v3++];
20                v5 = sub_21860(
21                    v1->ptr_field_28,
22                    v4->dat_field_10,
23                    sub_18F30, v4
24                );
25                v6 = v1->ptr_field_28;
26                v4->ptr_field_18 = v5;
27                sub_219C0(v6, v5, 1);
28            }
29        }
30        result = 0;
31    }
32    return result;
33 }
```

IDA Pro w/ our tool

Evaluation: Argument Decomposition Results

```

1 int network_rxxx(server *srv)
2 {
3     server *v1; // rbx
4     int result; // rax
5     size_t v3; // rbp
6     server_socket *v4; // r12
7     fdnode *v5; // rax
8     fdevents *v6; // rdi
9
10    v1 = srv;
11    result = fdevent_sxxx(srv->ev);
12    if ( result != -1 )
13    {
14        v3 = 0LL;
15        if ( !srv->sockets_disabled )
16        {
17            while ( v1->srv_sockets.
18                  used > v3 )
19            {
20                v4 = v1->srv_sockets.
21                      ptr[v3++];
22                v5 = fdevent_gxxx(
23                    v1->ev,
24                    v4->fd,
25                    network_sxxx, v4
26                );
27                v6 = v1->ev;
28                v4->fdn = v5;
29                fdevent_fxxx(v6, v5, 1);
30            }
31        }
32        result = 0LL;
33    }
34    return result;
35 }
```

IDA Pro w/ debug information

```

1 __int32 __fastcall sub_D840(__int64 a1)
2 {
3     _QWORD *v1; // rbx
4     __int32 result; // rax
5     unsigned __int64 v3; // rbp
6     int64 v4; // r12
7     __int64 v5; // rax
8     __int64 v6; // rdi
9
10    v1 = (_QWORD *)a1;
11    result = sub_12B7A(*(_QWORD *)(a1 + 24));
12    if ( (_DWORD)result != -1 )
13    {
14        v3 = 0LL;
15        if ( !*(__DWORD *)(a1 + 100) )
16        {
17            while ( v1[2] > v3 )
18            {
19                v4 = *(_QWORD *)(*v1 + 8 * v3++);
20                v5 = sub_21860(
21                    v1[3],
22                    *(unsigned int *)(v4 + 112),
23                    sub_18F30, v4
24                );
25                v6 = v1[3];
26                *(_QWORD *) (v4 + 120) = v5;
27                sub_219C0(v6, v5, 1);
28            }
29            result = 0;
30        }
31    }
32    return result;
33 }
```

Vanilla IDA Pro

```

1 __int32 __fastcall sub_D840(struct_C264 *a1)
2 {
3     struct_C264 *v1; // rbx
4     __int32 result; // rax
5     unsigned __int64 v3; // rbp
6     struct_CF4A *v4; // r12
7     struct_12A42 *v5; // rax
8     struct_12A0E *v6; // rdi
9
10    v1 = a1;
11    result = sub_12B7A(a1->ptr_field_28);
12    if ( result != -1 )
13    {
14        v3 = 0LL;
15        if ( !a1->dat_field_74 )
16        {
17            while ( v1->dat_field_10 > v3 )
18            {
19                v4 = v1->ptr_ptr_field_0[v3++];
20                v5 = sub_21860(
21                    v1->ptr_field_28,
22                    v4->dat_field_10,
23                    sub_18F30, v4
24                );
25                v6 = v1->ptr_field_28;
26                v4->ptr_field_18 = v5;
27                sub_219C0(v6, v5, 1);
28            }
29        }
30        result = 0;
31    }
32    return result;
33 }
```

IDA Pro w/ our tool

Evaluation: Argument Decomposition Results

```

1 int network_rxxx(server *srv)
2 {
3     server *v1; // rbx
4     int result; // rax
5     size_t v3; // rbp
6     server_socket *v4; // r12
7     fdnode *v5; // rax
8     fdevents *v6; // rdi
9
10    v1 = srv;
11    result = fdevent_sxxx(srv->ev);
12    if ( result != -1 )
13    {
14        v3 = 0LL;
15        if ( !srv->sockets_disabled )
16        {
17            while ( v1->srv_sockets.
18                  used > v3 )
19            {
20                v4 = v1->srv_sockets.
21                      ptr[v3++];
22                v5 = fdevent_gxxx(
23                    v1->ev,
24                    v4->fd,
25                    network_sxxx, v4
26                );
27                v6 = v1->ev;
28                v4->fdn = v5;
29                fdevent_fxxx(v6, v5, 1);
30            }
31        }
32        result = 0LL;
33    }
34    return result;
35 }
```

IDA Pro w/ debug information

```

1 __int32 __fastcall sub_D840(__int64 a1)
2 {
3     _QWORD *v1; // rbx
4     __int32 result; // rax
5     unsigned __int64 v3; // rbp
6     __int64 v4; // r12
7     __int64 v5; // rax
8     __int64 v6; // rdi
9
10    v1 = (_QWORD *)a1;
11    result = sub_12B7A(*(_QWORD *)(a1 + 24));
12    if ( (_DWORD)result != -1 )
13    {
14        v3 = 0LL;
15        if ( !*(__DWORD *)(a1 + 100) )
16        {
17            while ( v1[2] > v3 )
18            {
19                v4 = *(_QWORD *)(*v1 + 8 * v3++);
20                v5 = sub_21860(
21                    v1[3],
22                    *(unsigned int *)(v4 + 112),
23                    sub_18F30, v4
24                );
25                v6 = v1[3];
26                *(_QWORD *) (v4 + 120) = v5;
27                sub_219C0(v6, v5, 1);
28            }
29        }
30        result = 0;
31    }
32    return result;
33 }
```

Vanilla IDA Pro

```

1 __int32 __fastcall sub_D840(struct_C264 *a1)
2 {
3     struct_C264 *v1; // rbx
4     __int32 result; // rax
5     unsigned __int64 v3; // rbp
6     struct_CF4A *v4; // r12
7     struct_12A42 *v5; // rax
8     struct_12A0E *v6; // rdi
9
10    v1 = a1;
11    result = sub_12B7A(a1->ptr_field_28);
12    if ( result != -1 )
13    {
14        v3 = 0LL;
15        if ( !a1->dat_field_74 )
16        {
17            while ( v1->dat_field_10 > v3 )
18            {
19                v4 = v1->ptr_ptr_field_0[v3++];
20                v5 = sub_21860(
21                    v1->ptr_field_28,
22                    v4->dat_field_10,
23                    sub_18F30, v4
24                );
25                v6 = v1->ptr_field_28;
26                v4->ptr_field_18 = v5;
27                sub_219C0(v6, v5, 1);
28            }
29        }
30        result = 0;
31    }
32    return result;
33 }
```

IDA Pro w/ our tool

Evaluation: Argument Decomposition Results

```

1 int network_rxxx(server *srv)
2 {
3     server *v1; // rbx
4     int result; // rax
5     size_t v3; // rbp
6     server_socket *v4; // r12
7     fdnode *v5; // rax
8     fdevents *v6; // rdi
9
10    v1 = srv;
11    result = fdevent_sxxx(srv->ev);
12    if ( result != -1 )
13    {
14        v3 = 0LL;
15        if ( !srv->sockets_disabled )
16        {
17            while ( v1->srv_sockets.
18                  used > v3 )
19            {
20                v4 = v1->srv_sockets.
21                      ptr[v3++];
22                v5 = fdevent_gxxx(
23                    v1->ev,
24                    v4->fd,
25                    network_sxxx, v4
26                );
27                v6 = v1->ev;
28                v4->fdn = v5;
29                fdevent_fxxx(v6, v5, 1);
30            }
31        }
32        result = 0LL;
33    }
34    return result;
35 }
```

IDA Pro w/ debug information

```

1 __int32 __fastcall sub_D840(__int64 a1)
2 {
3     _QWORD *v1; // rbx
4     __int32 result; // rax
5     unsigned __int64 v3; // rbp
6     __int64 v4; // r12
7     __int64 v5; // rax
8     __int64 v6; // rdi
9
10    v1 = (_QWORD *)a1;
11    result = sub_12B7A(*(_QWORD *)(a1 + 24));
12    if ( (_DWORD)result != -1 )
13    {
14        v3 = 0LL;
15        if ( !*(__DWORD *)(a1 + 100) )
16        {
17            while ( v1[2] > v3 )
18            {
19                v4 = *(_QWORD *)(*v1 + 8 * v3++);
20                v5 = sub_21860(
21                    v1[3],
22                    *(unsigned int *)(v4 + 112),
23                    sub_18F30, v4
24                );
25                v6 = v1[3];
26                *(_QWORD *) (v4 + 120) = v5;
27                sub_219C0(v6, v5, 1);
28            }
29            result = 0;
30        }
31    }
32    return result;
33 }
```

Vanilla IDA Pro

```

1 __int32 __fastcall sub_D840(struct_C264 *a1)
2 {
3     struct_C264 *v1; // rbx
4     __int32 result; // rax
5     unsigned __int64 v3; // rbp
6     struct_CF4A *v4; // r12
7     struct_12A42 *v5; // rax
8     struct_12A0E *v6; // rdi
9
10    v1 = a1;
11    result = sub_12B7A(a1->ptr_field_28);
12    if ( result != -1 )
13    {
14        v3 = 0LL;
15        if ( !a1->dat_field_74 )
16        {
17            while ( v1->dat_field_10 > v3 )
18            {
19                v4 = v1->ptr_ptr_field_0[v3++];
20                v5 = sub_21860(
21                    v1->ptr_field_28,
22                    v4->dat_field_10,
23                    sub_18F30, v4
24                );
25                v6 = v1->ptr_field_28;
26                v4->ptr_field_18 = v5;
27                sub_219C0(v6, v5, 1);
28            }
29        }
30        result = 0;
31    }
32    return result;
33 }
```

IDA Pro w/ our tool

Evaluation: Argument Decomposition Results

```

1 int network_rxxx(server *srv)
2 {
3     server *v1; // rbx
4     int result; // rax
5     size_t v3; // rbp
6     server_socket *v4; // r12
7     fdnode *v5; // rax
8     fdevents *v6; // rdi
9
10    v1 = srv;
11    result = fdevent_sxxx(srv->ev);
12    if ( result != -1 )
13    {
14        v3 = 0LL;
15        if ( !srv->sockets_disabled )
16        {
17            while ( v1->srv_sockets.
18                   used > v3 )
19            {
20                v4 = v1->srv_sockets.
21                      ptr[v3++];
22                v5 = fdevent_gxxx(
23                    v1->ev,
24                    v4->fd,
25                    network_sxxx, v4
26                );
27                v6 = v1->ev;
28                v4->fdn = v5;
29                fdevent_fxxx(v6, v5, 1);
30            }
31        }
32        result = 0LL;
33    }
34    return result;
35 }
```

IDA Pro w/ debug information

```

1 __int32 __fastcall sub_D840(__int64 a1)
2 {
3     _QWORD *v1; // rbx
4     __int32 result; // rax
5     unsigned __int64 v3; // rbp
6     __int64 v4; // r12
7     __int64 v5; // rax
8     __int64 v6; // rdi
9
10    v1 = (_QWORD *)a1;
11    result = sub_12B7A(*(_QWORD *)(a1 + 24));
12    if ( (_DWORD)result != -1 )
13    {
14        v3 = 0LL;
15        if ( !*(__DWORD *)(a1 + 100) )
16        {
17            while ( v1[2] > v3 )
18            {
19                v4 = *(_QWORD *)(*v1 + 8 * v3++);
20
21                v5 = sub_21860(
22                    v1[3],
23                    *(unsigned int *) (v4 + 112),
24                    sub_18F30, v4
25                );
26                v6 = v1[3];
27                *(_QWORD *) (v4 + 120) = v5;
28                sub_219C0(v6, v5, 1);
29            }
30        }
31        result = 0;
32    }
33    return result;
34 }
```

Vanilla IDA Pro

```

1 __int32 __fastcall sub_D840(struct_C264 *a1)
2 {
3     struct_C264 *v1; // rbx
4     __int32 result; // rax
5     unsigned __int64 v3; // rbp
6     struct_CF4A *v4; // r12
7     struct_12A42 *v5; // rax
8     struct_12A0E *v6; // rdi
9
10    v1 = a1;
11    result = sub_12B7A(a1->ptr_field_28);
12    if ( result != -1 )
13    {
14        v3 = 0LL;
15        if ( !a1->dat_field_74 )
16        {
17            while ( v1->dat_field_10 > v3 )
18            {
19                v4 = v1->ptr_ptr_field_0[v3++];
20
21                v5 = sub_21860(
22                    v1->ptr_field_28,
23                    v4->dat_field_10,
24                    sub_18F30, v4
25                );
26                v6 = v1->ptr_field_28;
27                v4->ptr_field_18 = v5;
28                sub_219C0(v6, v5, 1);
29            }
30        }
31        result = 0;
32    }
33    return result;
34 }
```

IDA Pro w/ our tool

Evaluation: Argument Decomposition Results

```

1 int network_rxxx(server *srv)
2 {
3     server *v1; // rbx
4     int result; // rax
5     size_t v3; // rbp
6     server_socket *v4; // r12
7     fnode *v5; // rax
8     fdevents *v6; // rdi
9
10    v1 = srv;
11    result = fdevent_sxxx(srv->ev);
12    if ( result != -1 )
13    {
14        v3 = 0LL;
15        if ( !srv->sockets_disabled )
16        {
17            while ( v1->srv_sockets.
18                  used > v3 )
19            {
20                v4 = v1->srv_sockets.
21                      ptr[v3++];
22                v5 = fdevent_gxxx(
23                    v1->ev,
24                    v4->fd,
25                    network_sxxx, v4
26                );
27                v6 = v1->ev;
28                v4->fdn = v5;
29                fdevent_fxxx(v6, v5, 1);
30            }
31        }
32        result = 0LL;
33    }
34    return result;
35 }
```

IDA Pro w/ debug information

```

1 __int32 __fastcall sub_D840(__int64 a1)
2 {
3     _QWORD *v1; // rbx
4     __int32 result; // rax
5     unsigned __int64 v3; // rbp
6     __int64 v4; // r12
7     __int64 v5; // rax
8     __int64 v6; // rdi
9
10    v1 = (_QWORD *)a1;
11    result = sub_12B7A(*(_QWORD *)(a1 + 24));
12    if ( (_DWORD)result != -1 )
13    {
14        v3 = 0LL;
15        if ( !*(__DWORD *)(a1 + 100) )
16        {
17            while ( v1[2] > v3 )
18            {
19                v4 = *(_QWORD *)(*v1 + 8 * v3++);
20                v5 = sub_21860(
21                    v1[3],
22                    *(unsigned int *) (v4 + 112),
23                    sub_18F30, v4
24                );
25                v6 = v1[3];
26                *(_QWORD *) (v4 + 120) = v5;
27                sub_219C0(v6, v5, 1);
28            }
29            result = 0;
30        }
31    }
32    return result;
33 }
```

Vanilla IDA Pro

```

1 __int32 __fastcall sub_D840(struct_C264 *a1)
2 {
3     struct_C264 *v1; // rbx
4     __int32 result; // rax
5     unsigned __int64 v3; // rbp
6     struct_CF4A *v4; // r12
7     struct_12A42 *v5; // rax
8     struct_12A0E *v6; // rdi
9
10    v1 = a1;
11    result = sub_12B7A(a1->ptr_field_28);
12    if ( result != -1 )
13    {
14        v3 = 0LL;
15        if ( !a1->dat_field_74 )
16        {
17            while ( v1->dat_field_10 > v3 )
18            {
19                v4 = v1->ptr_ptr_field_0[v3++];
20                v5 = sub_21860(
21                    v1->ptr_field_28,
22                    v4->dat_field_10,
23                    sub_18F30, v4
24                );
25                v6 = v1->ptr_field_28;
26                v4->ptr_field_18 = v5;
27                sub_219C0(v6, v5, 1);
28            }
29        }
30        result = 0;
31    }
32    return result;
33 }
```

IDA Pro w/ our tool

Related Works

Variable recovery and type inference for stripped binary:

- Slowinska, Asia, Traian Stancescu, and Herbert Bos. "Howard: A Dynamic Excavator for Reverse Engineering Data Structures." *NDSS*. 2011.
- Shoshtaishvili, Yan, et al. "Sok:(state of) the art of war: Offensive techniques in binary analysis." *2016 IEEE Symposium on Security and Privacy (SP)*. IEEE, 2016.
- Lee, JongHyup, Thanassis Avgerinos, and David Brumley. "TIE: Principled reverse engineering of types in binary programs." (2011).
- Lin, Zhiqiang, Xiangyu Zhang, and Dongyan Xu. "Automatic reverse engineering of data structures from binary execution." *Proceedings of the 11th Annual Information Security Symposium*. 2010.

Probabilistic Program Analysis:

- Geldenhuis, Jaco, Matthew B. Dwyer, and Willem Visser. "Probabilistic symbolic execution." *Proceedings of the 2012 International Symposium on Software Testing and Analysis*. 2012.
- Borges, Mateus, et al. "Iterative distribution-aware sampling for probabilistic symbolic execution." *Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering*. 2015.
- Kwiatkowska, Marta, Gethin Norman, and David Parker. "PRISM 4.0: Verification of probabilistic real-time systems." *International conference on computer aided verification*. Springer, Berlin, Heidelberg, 2011.
- Filieri, Antonio, Carlo Ghezzi, and Giordano Tamburrelli. "Run-time efficient probabilistic model checking." *2011 33rd International Conference on Software Engineering (ICSE)*. IEEE, 2011.

We develop a novel probabilistic variable and data structure recovery technique for stripped binaries.

- It leverages an advanced data-flow technique, BDA, to collect program behaviors.
- It features using random variables to denote the likelihood of recovery results such that a large number of various kinds of hints can be integrated with the inherent uncertainty considered.
- It achieves the state-of-the-art inversing results.

We develop a novel probabilistic variable and data structure recovery technique for stripped binaries.

- It leverages an advanced data-flow technique, BDA, to collect program behaviors.
- It features using random variables to denote the likelihood of recovery results such that a large number of various kinds of hints can be integrated with the inherent uncertainty considered.
- It achieves the state-of-the-art inversing results.

Thanks!