

# *Constraint-Guided Statistical Type Reconstruction for Decompilation*

*Jeremy Lacomis*  
Carnegie Mellon University





## 1. jlacomis@gs17931:~/Data/coreutils/debug/src (ssh)

40299c:	89 f0	mov	%esi,%eax
40299e:	83 e0 04	and	\$0x4,%eax
4029a1:	74 1d	je	4029c0 <main+0x8b0>
4029a3:	31 d2	xor	%edx,%edx
4029a5:	48 89 d8	mov	%rbx,%rax
4029a8:	48 f7 f7	div	%rdi
4029ab:	48 89 05 be b8 20 00	mov	%rax,0x20b8be(%rip)
4029b2:	48 89 15 ff ba 20 00	mov	%rdx,0x20baff(%rip)
4029b9:	4d 85 c0	test	%r8,%r8
4029bc:	75 14	jne	4029d2 <main+0x8c2>
4029be:	eb 31	jmp	4029f1 <main+0x8e1>
4029c0:	48 83 fb ff	cmp	\$0xfffffffffffffff,%rbx
4029c4:	74 07	je	4029cd <main+0x8bd>
4029c6:	48 89 1d a3 b8 20 00	mov	%rbx,0x20b8a3(%rip)
4029cd:	4d 85 c0	test	%r8,%r8
4029d0:	74 1f	je	4029f1 <main+0x8e1>
4029d2:	89 c8	mov	%ecx,%eax
4029d4:	83 e0 10	and	\$0x10,%eax
4029d7:	74 18	je	4029f1 <main+0x8e1>

# Reverse Engineering

26 Feb 2013 | 14:00 GMT

## The Real Story of Stuxnet

How Kaspersky Lab tracked down the malware that stymied Iran's nuclear-fuel enrichment program

By David Kushner



**Computer cables snake across the floor. Cryptic flowcharts are scrawled across various whiteboards adorning the walls. A life-size Batman doll stands in the hall. This office might seem no different than any other geeky workplace, but in fact it's the front line of a war—a cyberwar, where most battles play out not in remote jungles or deserts but in suburban office parks like this one.**

2014 IEEE International Conference on Software Maintenance and Evolution

## Reverse Engineering PL/SQL Legacy Code: An Experience Report

Martin Habringer  
voestalpine Stahl GmbH  
4020 Linz, Austria  
martin.habringer@voestalpine.com

Michael Moser and Josef Pichler  
Software Analytics and Evolution  
Software Competence Center Hagenberg GmbH  
4232 Hagenberg, Austria  
michael.moser@scch.at, josef.pichler@scch.at

- Changes in business cases over the last years were not reflected in verification logic of the legacy code.
- For a new production plant, additional requirements must be incorporated.
- The maintenance of the legacy programs was complicated by the retirement of original developers.
- Legacy code is not extensible in a safe and reliable way.
- Stakeholders estimated high effort for manual analysis of the legacy code.

The goal for the reverse engineering tool was to support stakeholders to comprehend the verification logic implemented in the legacy programs. Whereas, comprehension requires that stakeholders can (1) *identify* the business cases currently checked by the software as well as that stakeholders are able to (2) *extend* the verification logic with respect to new requirements.

The contributions of this paper are:

- 1) Presenting requirements for a reverse engineering tool aimed at PL/SQL legacy code.

IDA - dd /media/DATA/coreutils/debug/src/dd

File Edit Jump Search View Debugger Options Windows Help

Library function Regular function Instruction Data Unexplored External symbol

Graph overview IDA View-A

xor edx, edx  
div rdi  
mov cs:skip\_records, eax  
mov cs:skip\_bytes, edx  
cmp rbx, 0xFFFFFFFFFFFFFFFh  
jnz short loc\_40299C

loc\_402996:  
cmp rbx, 0xFFFFFFFFFFFFFFFh  
jz short loc\_4029C0

jmp short loc\_4029C0

loc\_40299C:  
mov eax, esi  
and eax, 4  
jz short loc\_4029C0

loc\_4029C0:  
cmp rbx, 0xFFFFFFFFFFFFFFFh  
jz short loc\_4029CD

100.00% (3865,13067) (657,9) 00002996 0000000000402996: main:loc\_402996

The initial autoanalysis has been finished.

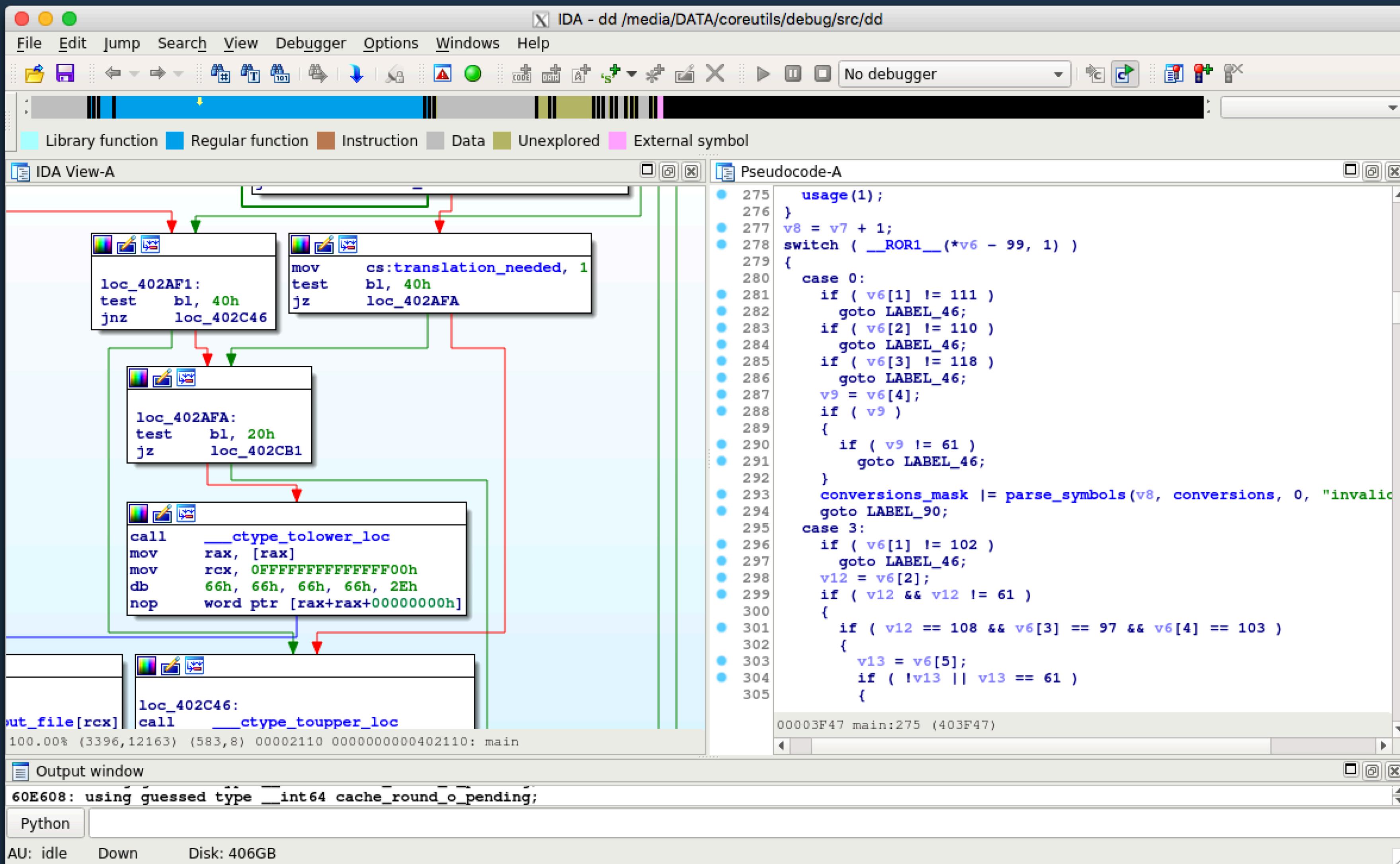
Output window

Python

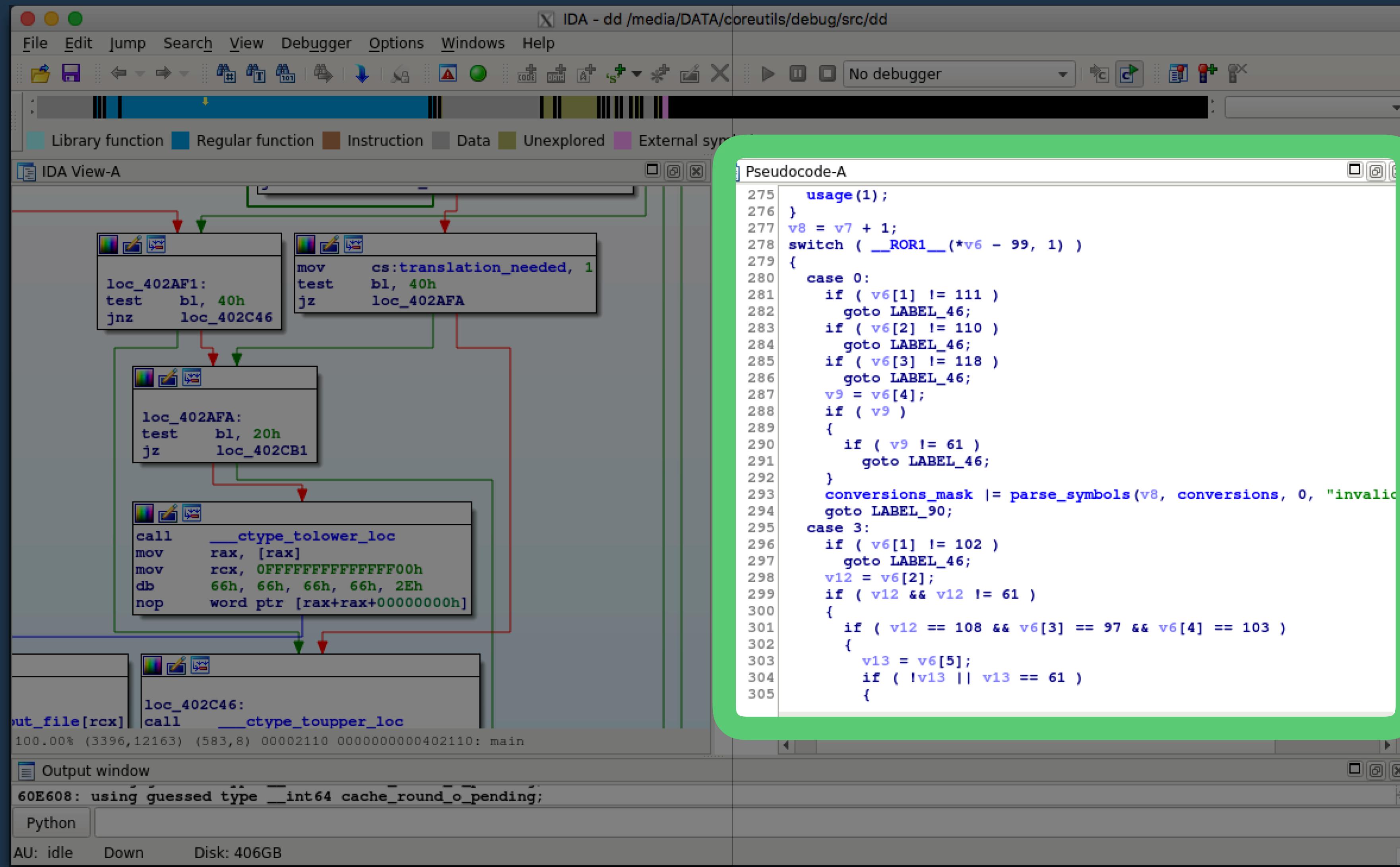
AU: idle Down Disk: 406GB

This screenshot shows the IDA Pro interface during a static analysis session. The main window displays a complex control flow graph (CFG) for the 'dd' program. The graph consists of several nodes, each containing assembly code. The nodes are interconnected by colored edges representing control flow: red for jumps, green for calls, and blue for returns. A large, semi-transparent red rectangle highlights a specific section of the graph, likely indicating a region of interest or a potential exploit path. The assembly code within these nodes includes instructions like XOR, DIV, MOV, CMP, and JNZ, along with labels such as loc\_402996, loc\_40299C, loc\_4029C0, and loc\_4029CD. The status bar at the bottom provides performance metrics and disk usage information.

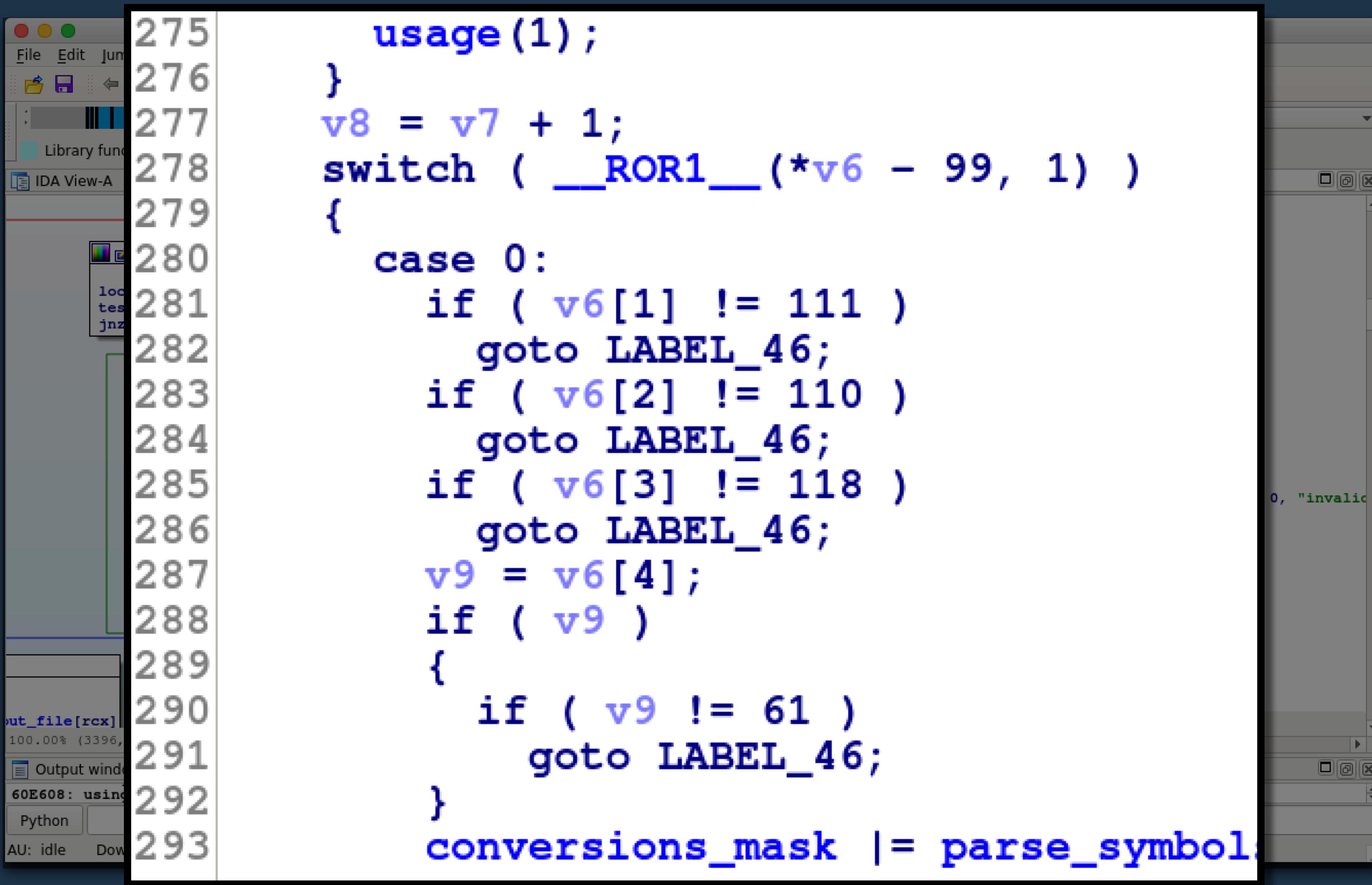
# Decompiler



# Decompiler



# Decompiler



The screenshot shows the IDA Pro decompiler window with the following assembly code:

```
275     usage(1);
276 }
277 v8 = v7 + 1;
278 switch ( __ROR1__(*v6 - 99, 1) )
279 {
280     case 0:
281         if ( v6[1] != 111 )
282             goto LABEL_46;
283         if ( v6[2] != 110 )
284             goto LABEL_46;
285         if ( v6[3] != 118 )
286             goto LABEL_46;
287         v9 = v6[4];
288         if ( v9 )
289         {
290             if ( v9 != 61 )
291                 goto LABEL_46;
292         }
293     conversions_mask |= parse_symbols;
```

The code is annotated with line numbers from 275 to 293. It includes several conditional jumps labeled `LABEL_46`. The assembly code is displayed in blue, while comments and labels are in black. The IDA Pro interface is visible around the code window.

# Decompilers don't Recover Code!

1. emacs-26.1

```
1#include <stdio.h>
2
3int main() {
4    int x = 0;
5    int y = 0;
6    while (x < 100) {
7        printf("%d\n", x);
8        x++;
9    }
10   return y;
11}
```

-UUU:----F1 count.c All (12,0) (C/\*l

Pseudocode-A

```
1int __cdecl main(int argc, const char **argv)
2{
3    signed int i; // [rsp+8h] [rbp-8h]
4
5    for ( i = 0; i < 100; ++i )
6        printf("%d\n", (unsigned int)i, envp);
7    return 0;
8}
```

00000F44 \_main:8 (100000F44)

# Compilation Loses Information

## Comments:

```
/* This is the functionality  
 * you're looking for! */
```

## Loop Constructs:

```
while (x < 100) {...}  
  
for (i = 0; i < 100; ++i) {...}
```

## Variable Names:

```
int width, length;  
double volume;  
char *user_id;
```

## User-Defined Types:

```
typedef struct {  
    int x;  
    int y;  
} point;
```

# Compilation Loses Information

## Comments:

```
/* This is the functionality  
 * you're looking for! */
```

## Loop Constructs:

```
while (x < 100) {...}  
  
for (i = 0; i < 100; ++i) {...}
```

## Variable Names:

```
int width, length;  
double volume;  
char *user_id;
```

## User-Defined Types:

```
typedef struct {  
    int x;  
    int y;  
} point;
```

# Compilation Loses Information

## Comments:

```
/* This is the functionality  
 * you're looking for! */
```

## Loop Constructs:

```
while (x < 100) {...}  
  
for (i = 0; i < 100; ++i) {...}
```

## Variable Names:

```
int width, length;  
double volume;  
char *user_id;
```

## User-Defined Types:

```
typedef struct {  
    int x;  
    int y;  
} point;
```

# Types Improve Variable Renaming

```
unsigned int width, total, attendance;  
size_t name_length;
```

**Meaningful Variable Names for Decompiled Code: A Machine Translation Approach,**  
A. Jaffe, J. Lacomis, E. J. Schwartz, C. Le Goues, and B. Vasilescu, in ICPC, 2018

# Types Can Guide Understanding

```
double func(int *a1, int *a2) {  
    double v1, v2;  
  
    v1 = pow(*a1 - *a2, 2);  
    v2 = pow(a1[1] - a2[1]), 2);  
  
    return sqrt(v1 + v2);  
}
```

# Types Can Guide Understanding

```
double func(int *a1, int *a2) {  
    double v1, v2;  
  
    v1 = pow(*a1 - *a2, 2);  
    v2 = pow(a1[1] - a2[1]), 2);  
  
    return sqrt(v1 + v2);  
}
```

int pointer arguments

# Types Can Guide Understanding

```
double func(int *a1, int *a2) {  
    double v1, v2;  
  
    v1 = pow(*a1 - *a2);  
    v2 = pow((a1[1] - a2[1]), 2);  
  
    return sqrt(v1 + v2);  
}
```

Dereference and do some math

# Types Can Guide Understanding

```
double func(int *a1, int *a2) {  
    double v1, v2;  
  
    v1 = pow((*a1 - *a2), 2);  
    v2 = pow((a1[1] - a2[1]), 2);  
  
    return sqrt(v1 + v2);  
}
```

Index like an array?!?

# Types Can Guide Understanding

```
typedef struct {
    int x;
    int y;
} point;

double func(point *a1, point *a2) {
    double v1, v2;

    v1 = pow((a1->x - a2->x), 2);
    v2 = pow((a1->y - a2->y), 2);

    return sqrt(v1 + v2);
}
```

# Types Can Guide Understanding

```
double func(int *a1, int *a2) {  
    double v1, v2;  
  
    v1 = pow((*a1 - *a2), 2);  
    v2 = pow((a1[1] - a2[1]), 2);  
  
    return sqrt(v1 + v2);  
}
```

```
typedef struct {  
    int x;  
    int y;  
} point;
```

```
double func(point *a1, point *a2) {  
    double v1, v2;  
  
    v1 = pow((a1->x - a2->x), 2);  
    v2 = pow((a1->y - a2->y), 2);  
  
    return sqrt(v1 + v2);  
}
```

# Types Can Guide Understanding

```
double func(int *a1, int *a2) {  
    double v1, v2;  
  
    v1 = pow(*a1 - *a2), 2);  
    v2 = pow(a1[1] - a2[1]), 2);  
  
    return sqrt(v1 + v2);  
}
```

```
typedef struct {  
    int x;  
    int y;  
} point;  
  
double func(point *a1, point *a2) {  
    double v1, v2;  
  
    v1 = pow((a1->x - a2->x), 2);  
    v2 = pow((a1->y - a2->y), 2);  
  
    return sqrt(v1 + v2);  
}
```

# Types Can Guide Understanding

```
double func(int *a1, int *a2) {  
    double v1, v2;  
  
    v1 = pow((*a1 - *a2), 2);  
    v2 = pow((a1[1] - a2[1]), 2);  
  
    return sqrt(v1 + v2);  
}
```

```
typedef struct {  
    int x;  
    int y;  
} point;
```

```
double func(point *a1, point *a2) {  
    double v1, v2;  
  
    v1 = pow((a1->x - a2->x), 2);  
    v2 = pow((a1->y - a2->y), 2);  
  
    return sqrt(v1 + v2);  
}
```

# Types Can Guide Understanding

```
double func(int *a1, int *a2) {  
    double v1, v2;  
  
    v1 = pow(*a1 - *a2, 2);  
    v2 = pow(a1[1] - a2[1]), 2);  
  
    return sqrt(v1 + v2);  
}
```

```
typedef struct {  
    int x;  
    int y;  
} point;  
  
double func(point *a1, point *a2) {  
    double v1, v2;  
  
    v1 = pow((a1->x - a2->x), 2);  
    v2 = pow((a1->y - a2->y), 2);  
  
    return sqrt(v1 + v2);  
}
```

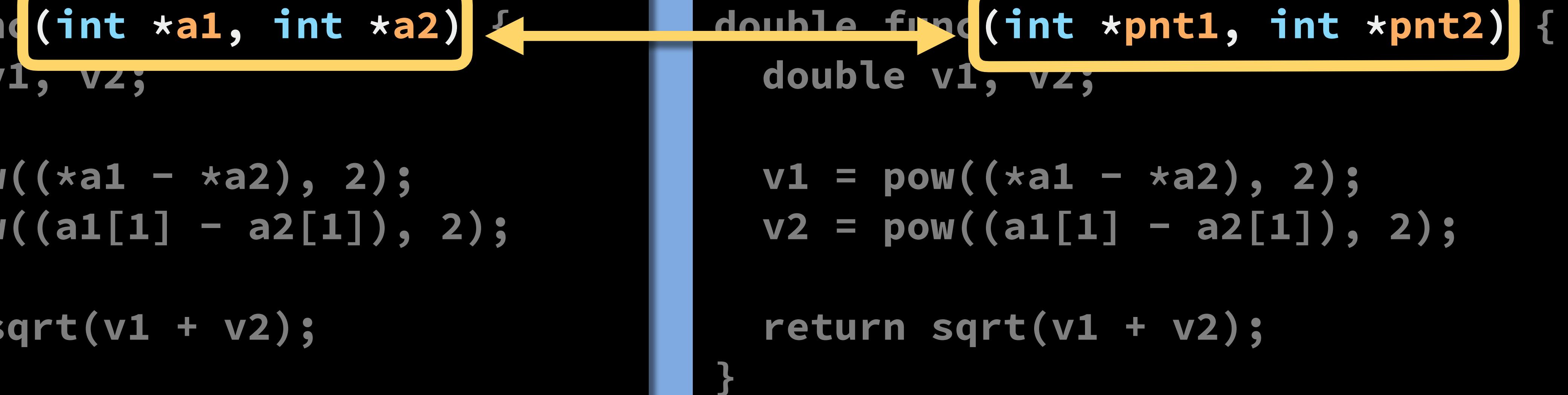
# Statistical: Select Types Using ML

```
double func(int *a1, int *a2) {  
    double v1, v2;  
  
    v1 = pow(*a1 - *a2, 2);  
    v2 = pow(a1[1] - a2[1]), 2);  
  
    return sqrt(v1 + v2);  
}
```

```
typedef struct {  
    int x;  
    int y;  
} point;  
  
double func(point *a1, point *a2) {  
    double v1, v2;  
  
    v1 = pow((a1->x - a2->x), 2);  
    v2 = pow((a1->y - a2->y), 2);  
  
    return sqrt(v1 + v2);  
}
```

# Statistical: Select Names Using ML

```
double func(int *a1, int *a2) {  
    double v1, v2;  
  
    v1 = pow(*a1 - *a2, 2);  
    v2 = pow(a1[1] - a2[1]), 2);  
  
    return sqrt(v1 + v2);  
}  
  
double func(int *pnt1, int *pnt2) {  
    double v1, v2;  
  
    v1 = pow(*pnt1 - *pnt2, 2);  
    v2 = pow(pnt1[1] - pnt2[1]), 2);  
  
    return sqrt(v1 + v2);  
}
```



# Constraint-Guided: Restrict to Compatible Types

```
double func(int *a1, int *a2) {  
    double v1, v2;  
  
    v1 = pow(*a1 - *a2, 2);  
    v2 = pow(a1[1] - a2[1]), 2);  
  
    return sqrt(v1 + v2);  
}
```

```
typedef struct {  
    int x;  
    int y;  
} point;
```

```
double func(point *a1, point *a2) {  
    double v1, v2;  
  
    v1 = pow((a1->x - a2->x), 2);  
    v2 = pow((a1->y - a2->y), 2);  
  
    return sqrt(v1 + v2);  
}
```

# Type Constraints

```
typedef struct {  
    int x;  
    int y;  
} point;  
  
sizeof(char);      // 1  
  
sizeof(int);       // 4  
sizeof(float);    // 4  
  
sizeof(point);     // 8  
sizeof(int[2]);   // 8
```

# Type Constraints

```
typedef struct {  
    int x;  
    int y;  
} point;  
  
sizeof(char); // 1  
  
sizeof(int); // 4  
sizeof(float); // 4  
  
sizeof(point); // 8  
sizeof(int[2]); // 8
```

# Type Constraints

```
typedef struct {  
    int x;  
    int y;  
} point;  
  
sizeof(char);      // 1  
  
sizeof(int);      // 4  
sizeof(float);    // 4  
  
sizeof(point);    // 8  
sizeof(int[2]);   // 8
```

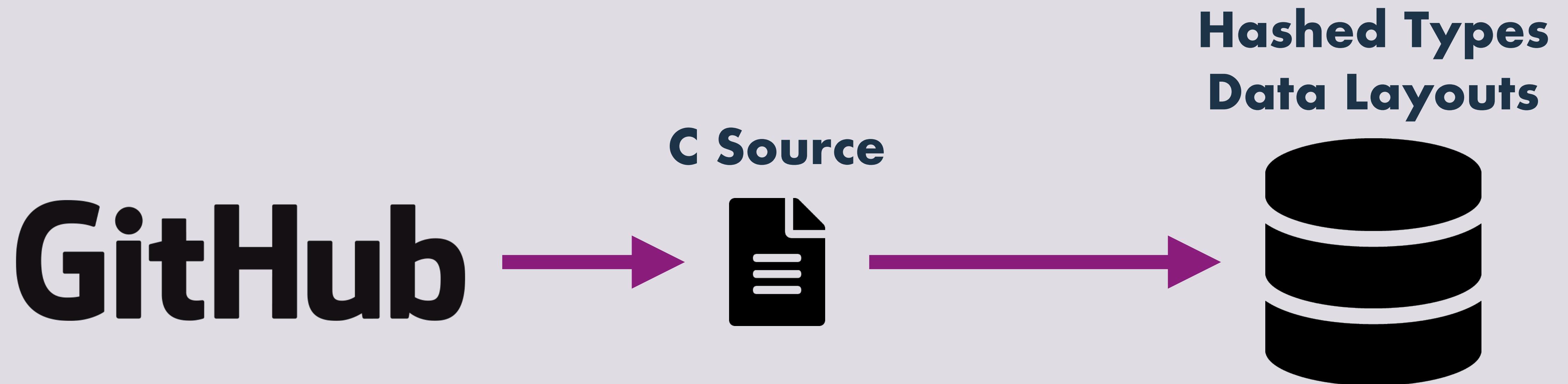
# Type Constraints

```
typedef struct {  
    int x;  
    int y;  
} point;  
  
sizeof(char);      // 1  
  
sizeof(int);       // 4  
sizeof(float);    // 4  
  
sizeof(point);     // 8  
sizeof(int[2]);   // 8
```

# Type Constraints

```
typedef struct {  
    int x;  
    int y;  
} point;  
  
sizeof(char);      // 1  
  
sizeof(int);       // 4  
sizeof(float);    // 4  
  
sizeof(point);     // 8  
sizeof(int[2]);   // 8
```

# Mining Types



**Meaningful Variable Names for Decompiled Code: A Machine Translation Approach,**  
A. Jaffe, J. Lacomis, E. J. Schwartz, C. Le Goues, and B. Vasilescu, in ICPC, 2018

# Statistical: Type Assignment using ML

```
double func(int *a1, int *a2) {  
    double v1, v2;  
  
    v1 = pow(*a1 - *a2, 2);  
    v2 = pow(a1[1] - a2[1]), 2);  
  
    return sqrt(v1 + v2);  
}
```

```
typedef struct {  
    int x;  
    int y;  
} point;  
  
double func(point *a1, point *a2) {  
    double v1, v2;  
  
    v1 = pow((a1->x - a2->x), 2);  
    v2 = pow((a1->y - a2->y), 2);  
  
    return sqrt(v1 + v2);  
}
```

# Previous Work: n-gram Models

I → am → very → tired  
I → am → very → banana

# Previous Work: n-gram Models

I → am → very → **tired**

I → am → very → **banana**

int → main → ( → int → **argc**

int → main → ( → int → **banana**

# Missing Context

## Decompiler Output

```
double func(int *a1, int *a2) {  
    double v1, v2;  
  
    v1 = pow((*a1 - *a2), 2);  
    v2 = pow((a1[1] - a2[1]), 2);  
  
    return sqrt(v1 + v2);  
}
```

## Original Code

```
double func(point *a1, point *a2) {  
    double v1, v2;  
  
    v1 = pow((a1->x - a2->x), 2);  
    v2 = pow((a1->y - a2->y), 2);  
  
    return sqrt(v1 + v2);  
}
```

# Missing Context

## Decompiler Output

```
double func(int *a1, int *a2) {  
    double v1, v2;  
  
    v1 = pow(*a1 - *a2, 2);  
    v2 = pow((a1[1] - a2[1]), 2);  
  
    return sqrt(v1 + v2);  
}
```

## Original Code

```
double func(point *a1, point *a2) {  
    double v1, v2;  
  
    v1 = pow((a1->x - a2->x), 2);  
    v2 = pow((a1->y - a2->y), 2);  
  
    return sqrt(v1 + v2);  
}
```

# Missing Context

## Decompiler Output

```
double func int *a1 int *a2 {
    double v1, v2;

    v1 = pow(*a1 - *a2, 2);
    v2 = pow(a1[1] - a2[1], 2);

    return sqrt(v1 + v2);
}
```

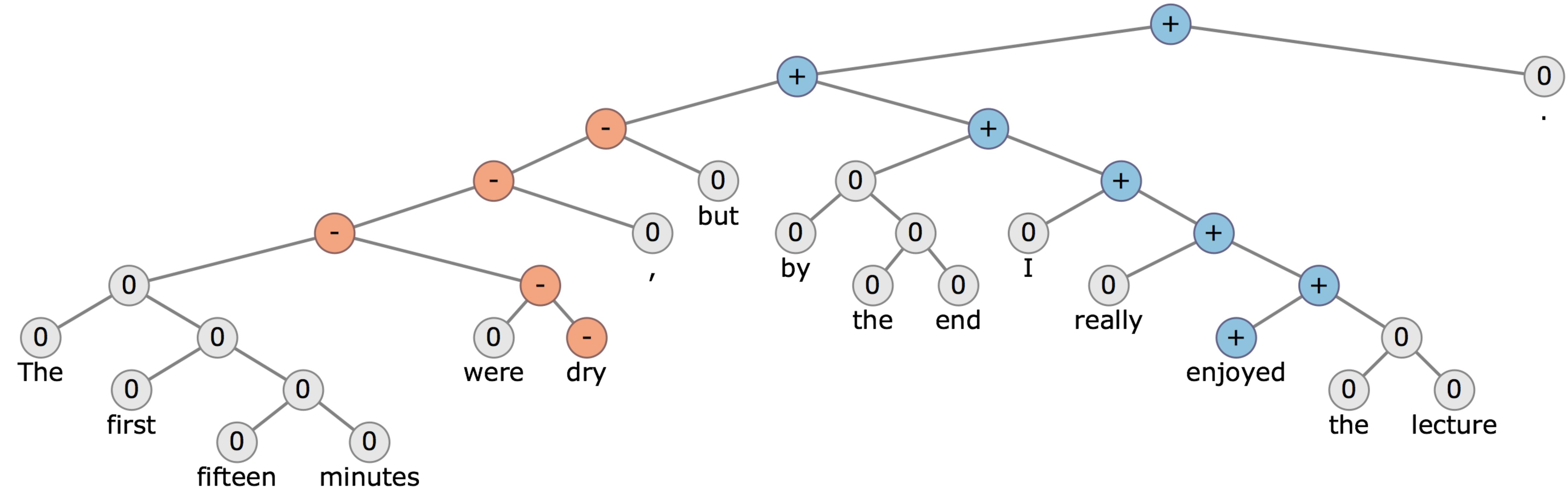
## Original Code

```
double func point *a1 point *a2 {
    double v1, v2;

    v1 = pow((a1->x - a2->x), 2);
    v2 = pow((a1->y - a2->y), 2);

    return sqrt(v1 + v2);
}
```

# Inspiration: Text Classification



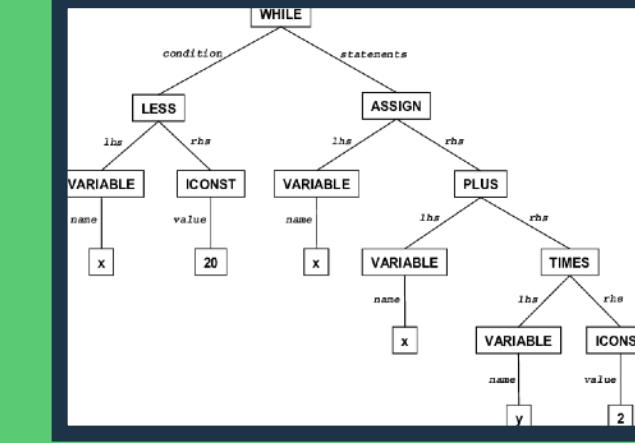
# Compilation

Source Code

```
double func(int *a1, int *a2) {...}
```

Parse

Abstract Syntax Tree



Generate Assembly

Assembly Code

```
i.jecomics@gs17931:~/Data/corutiles/debug/src (ssh)
46299c: 89 j4 09 04    mov    %esi,%eax
46299d: 89 e9 04    and    %eax,%eax
46299e: 74 1d        je     4629e0
46299f: 31 d2        xor    %edx,%edx
4629a0: 48 89 d8    mov    %rbx,%rax
4629a1: 48 f7 f7    div    %rdi
4629a2: 48 89 05 be b8 20 00  mov    %rax,%0xb0bbe(%rip)
4629a3: 48 89 15 ff ba 20 00  mov    %rdx,%0xb0bff(%rip)
4629a4: 4d c0        test   %rcx,%rcx
4629a5: 74 01        je     4629a9
4629a6: 48 31        cmp    %eax,%eax
4629a7: 48 83 fb ff    cmp    $0xfffffffffffffff,%rbx
4629a8: 74 07        je     4629cd
4629a9: 48 89 1d a3 b8 20 00  mov    %rbx,%0xb0ba3(%rip)
4629a9: 4d c0        test   %rcx,%rcx
4629a9: 74 1f        je     4629f1
4629a9: 89 c8        mov    %secx,%eax
4629a9: 83 e9 10    and    $0x10,%eax
4629a9: 74 18        je     4629f1
4629a9: 0f 1f 45 00    inc    %eax
```

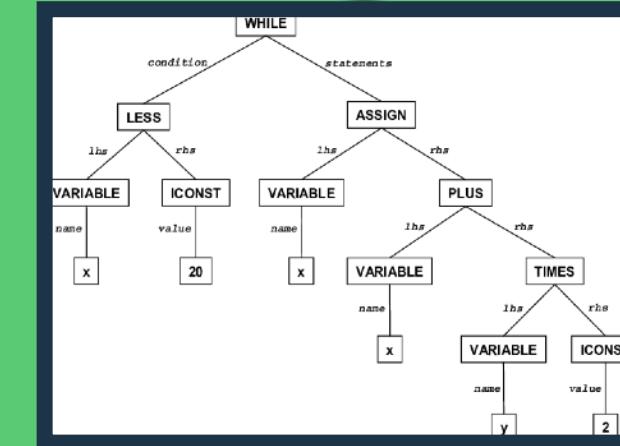
# Decompilation

Decompiled Code

```
double func(int *a1, int *a2) {...}
```

Generate Code

Abstract Syntax Tree



Generate Tree

Disassembled Code

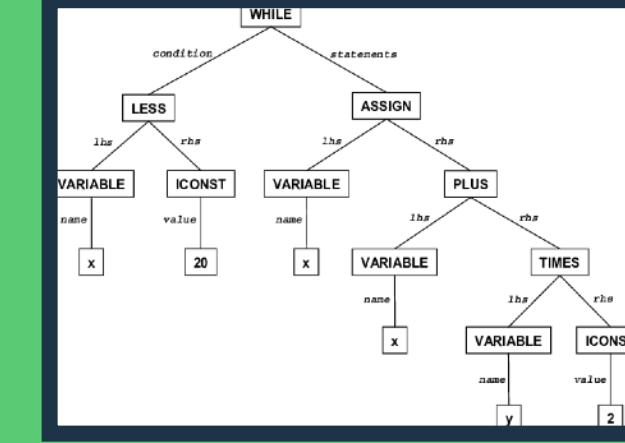
```
t.jecomics@gs17931:~/Data/coreutils/debug/src (ssh)
46299c: 89 e9 04 00 mov %esi,%eax
46299d: 89 e9 04 00 and %eax,%eax
4629a1: 74 1d je 4629e0
4629a3: 31 d2 xor %edx,%edx
4629a5: 48 89 d8 mov %rbx,%rax
4629a8: 48 f7 f7 div %rdi
4629ab: 48 89 05 be b8 20 00 mov %rax,%0xb0bbe(%rip)
4629b2: 48 89 15 ff ba 20 00 mov %rdx,%0xb0bff(%rip)
4629b9: 4d c0 test %eax,%eax
4629bc: 48 83 31 jne 4629c2
4629c0: 48 83 fb ff cmp $0xfffffffffffffff,%rbx
4629c4: 74 07 je 4629cd
4629c6: 48 89 1d a3 b8 20 00 mov %rbx,%0xb0ba3(%rip)
4629cd: 4d c0 test %rbx,%rbx
4629d0: 48 89 1f je 4629f1
4629d2: 89 c8 mov %ecx,%eax
4629d4: 83 e9 10 and $0x10,%eax
4629d7: 74 18 je 4629f1
```

# Decompilation

Decompiled Code

```
double func(int *a1, int *a2) {...}
```

Abstract Syntax Tree



Disassembled Code

```
l_jpcmini@gs17931:~/Data/corvella$ debugger/a.out
00299c: 89 f8          mov    %eax,%eax
00299d: 89 e6          mov    %eax,%ax
00299e: 74 04          je    0029a2<main+0x8b0>
0029a1: 31 d2          xor    %edx,%edx
0029a3: 31 d2          xor    %edx,%edx
0029a5: 48 89 d8        mov    %rax,%rdx
0029a8: 48 f7 f7        div    %rdi
0029ab: 48 89 05 be b8 20 00  mov    %rax,%rbx(%rip)
0029a2: 48 89 15 ff ba 20 00  mov    %rax,%rbx(%rip)
002999: 4d 85 c0        test   %rsi,%r8
00299c: 79 31          jne    0029a2<main+0x8b0>
00299d: 48 83 31          cmp    %rsi,%rsi
00299e: 48 83 fb ff        mov    %rbx,%rbx
0029a1: 74 07          je    0029cd<main+0x8bd>
0029a5: 48 89 1d a3 b8 20 00  mov    %rbx,%rbx(%rip)
0029a2: 4d 85 c0        test   %rsi,%r8
0029a5: 74 1f          je    0029f1<main+0x8e1>
0029a2: 89 c8          mov    %rcx,%eax
0029a3: 89 e0          and    %eax,%eax
0029a4: 74 18          je    0029f1<main+0x8e1>
```

Generate Code

Generate Tree

# Strategy

## Decompiler Output

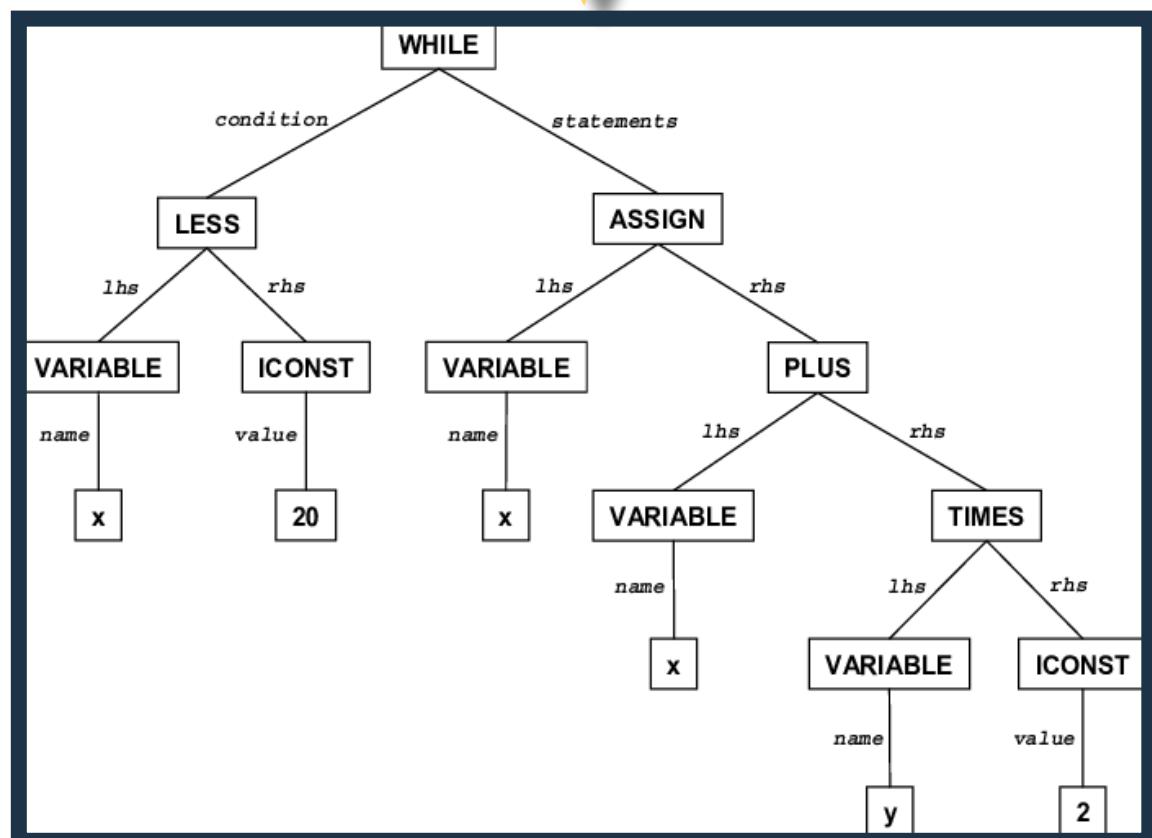
```
double func(int *a1, int *a2) {...}
```

# Strategy

## Decompiler Output

```
double func(int *a1, int *a2) {...}
```

Parse

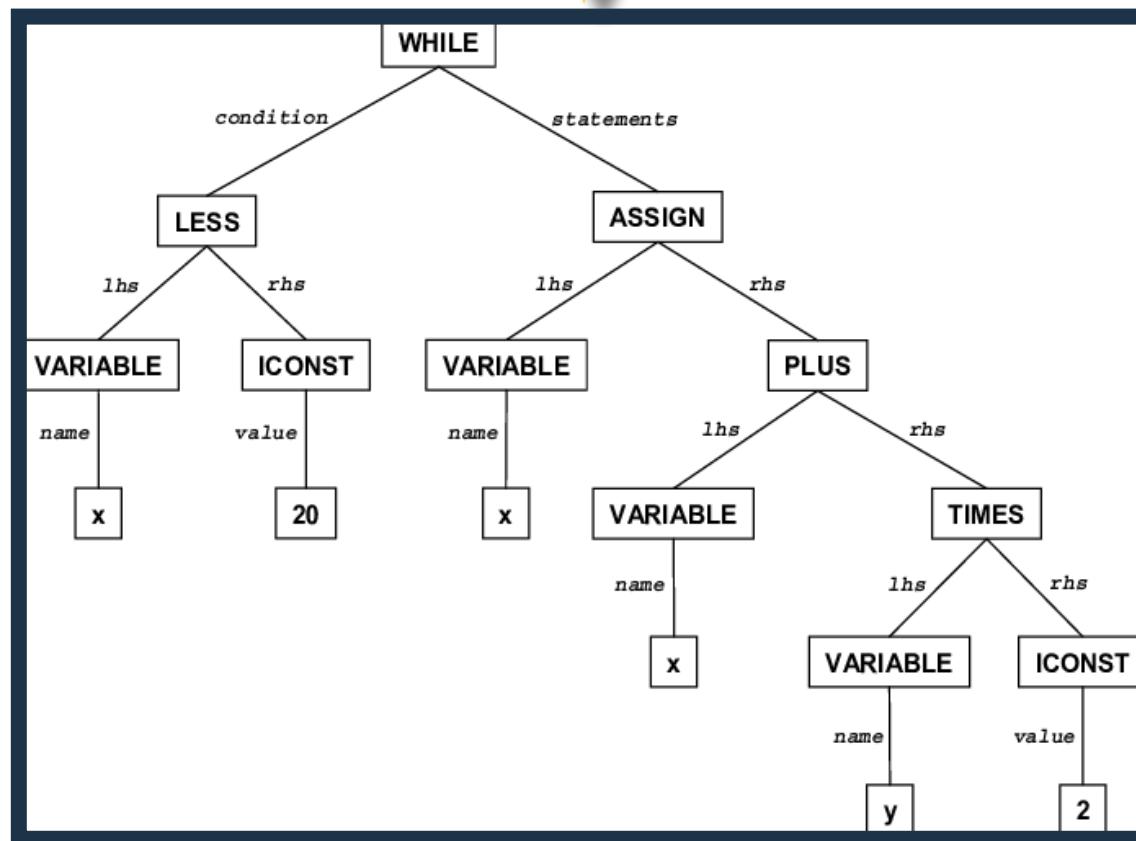


# Strategy

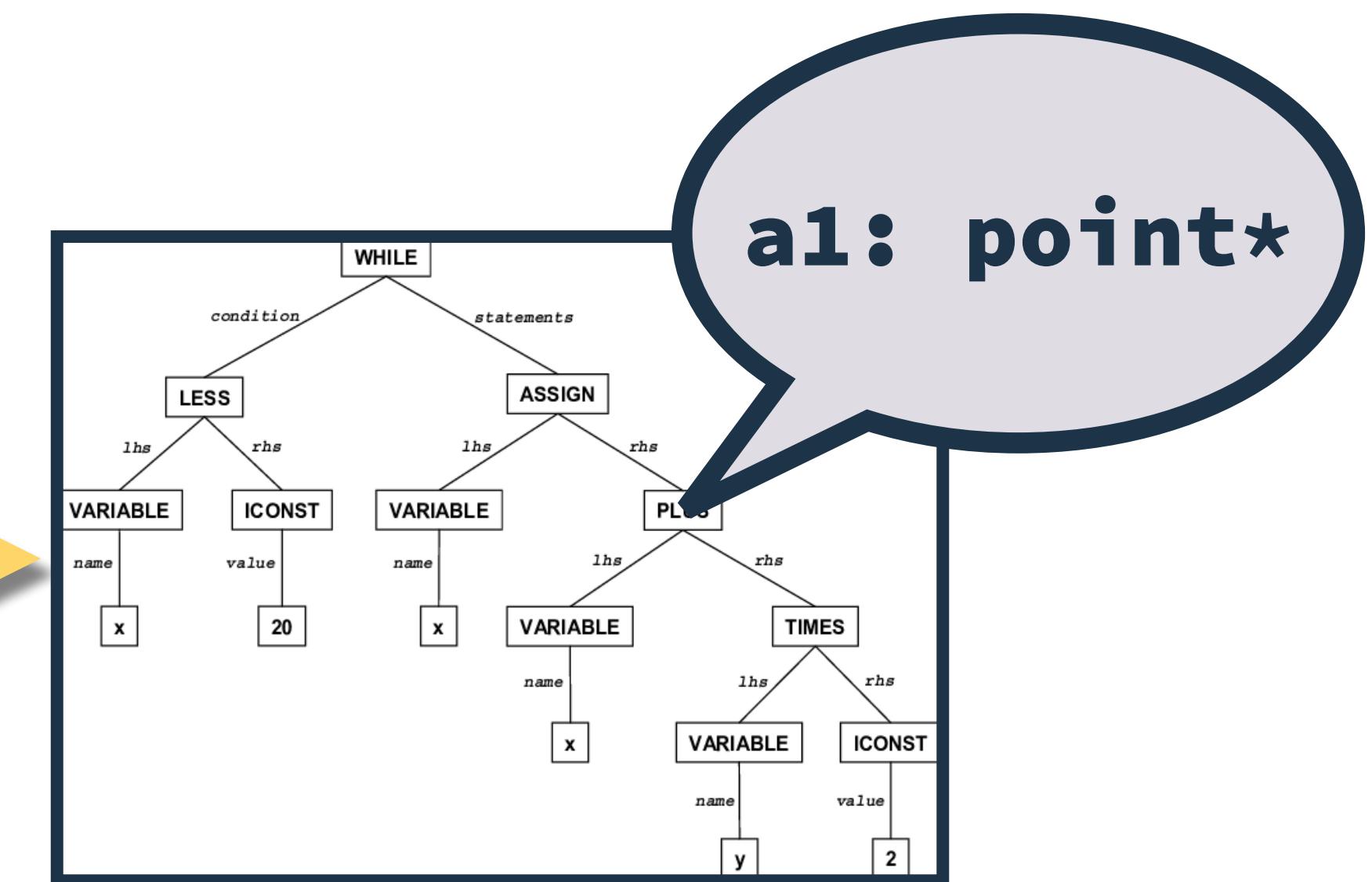
## Decompiler Output

```
double func(int *a1, int *a2) {...}
```

Parse



Tagger

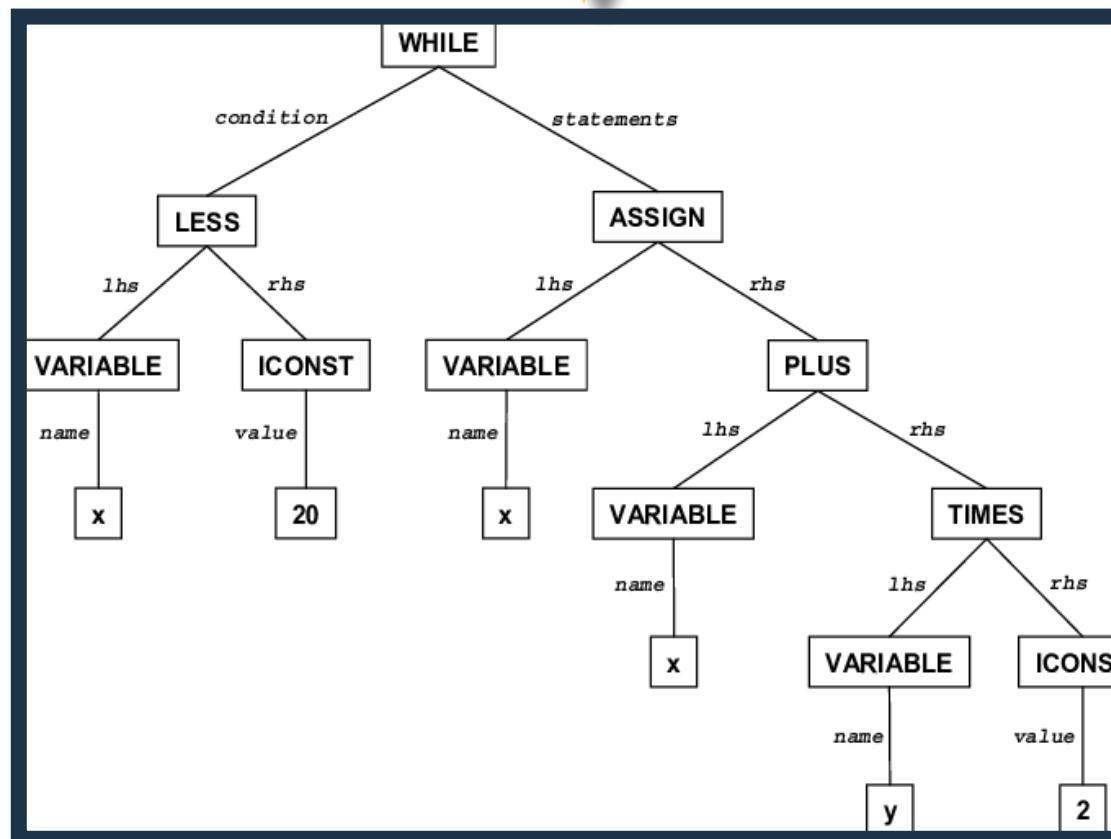


# Strategy

## Decompiler Output

```
double func(int *a1, int *a2) {...}
```

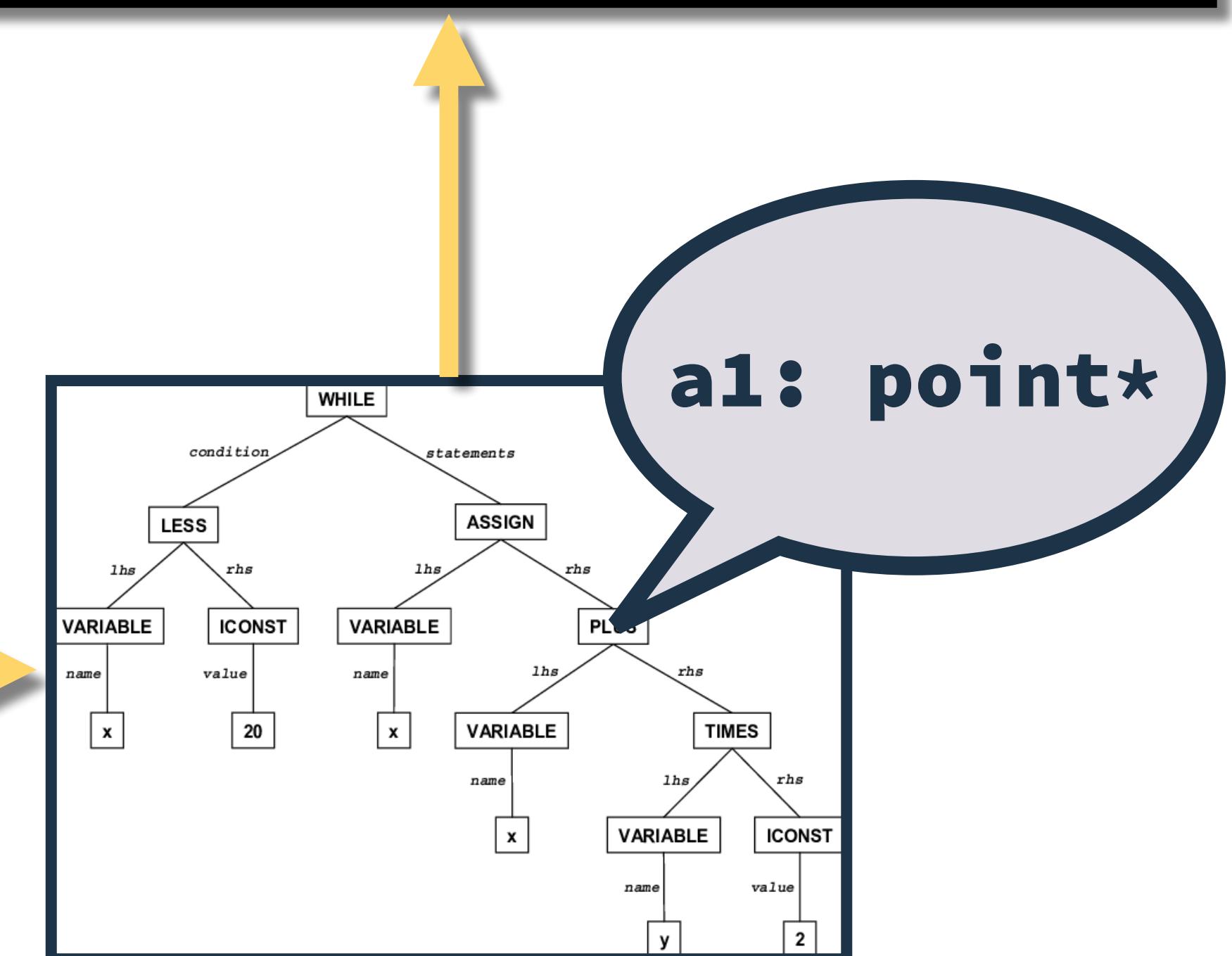
Parse



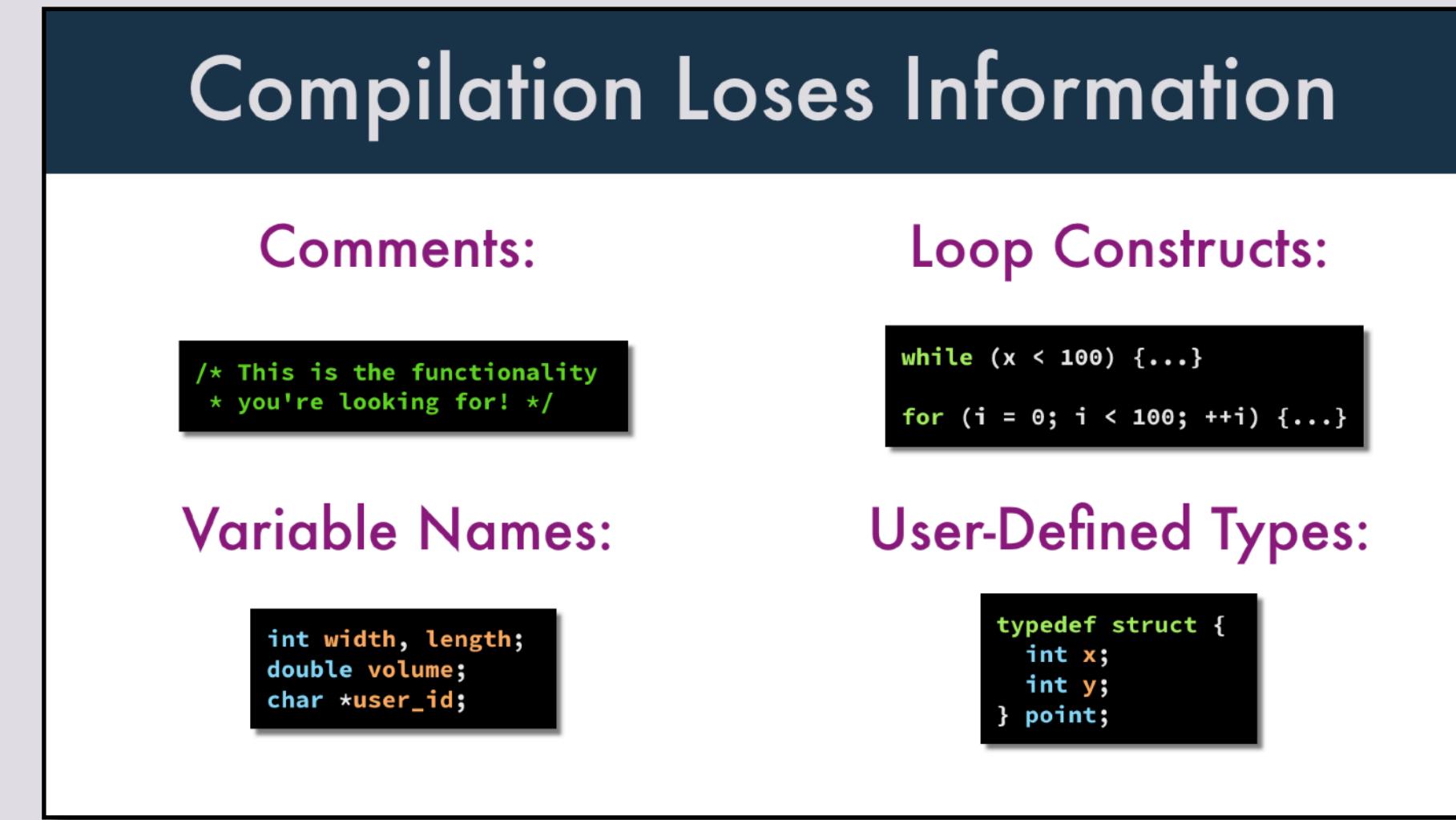
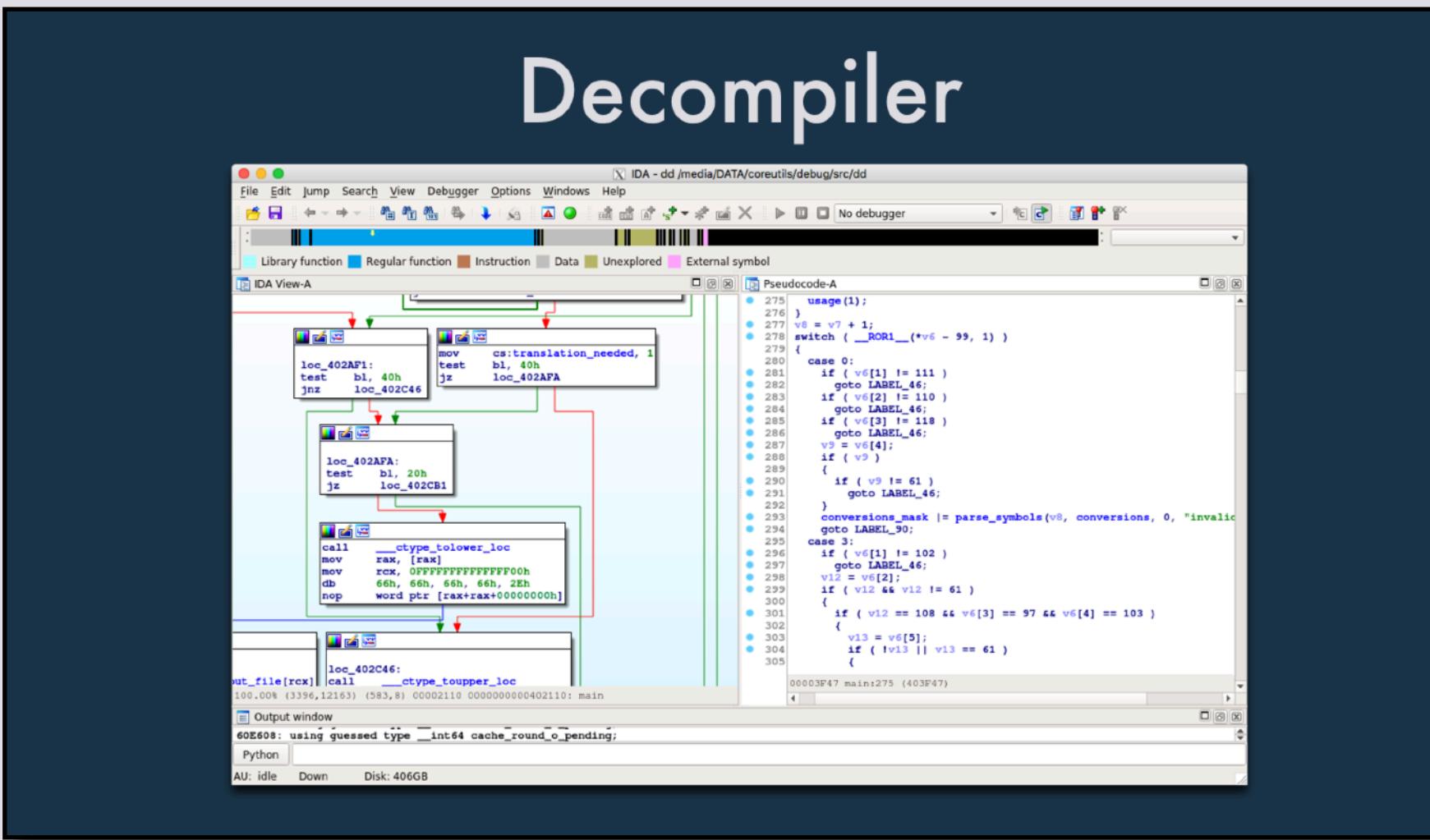
## New Output

```
double func(point *a1, int *a2) {...}
```

Tagger

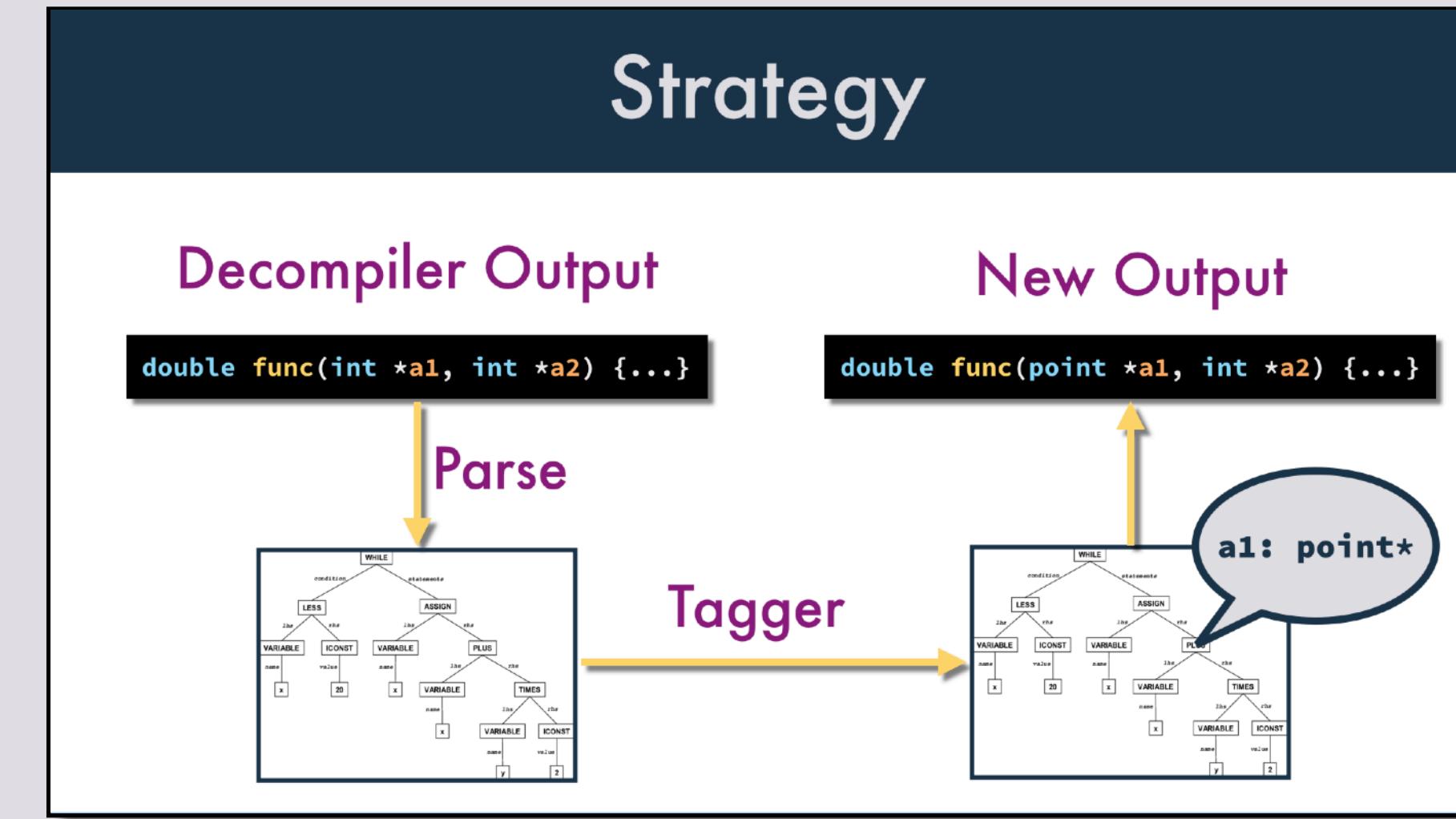


# Conclusion



## Type Constraints

```
typedef struct {  
    int x;  
    int y;  
} point;  
  
sizeof(char); // 1  
sizeof(int); // 4  
sizeof(float); // 4  
  
sizeof(point); // 8  
sizeof(int[2]); // 8
```



# Contact Me

**Jeremy Lacomis**



✉ [jlacomis@cmu.edu](mailto:jlacomis@cmu.edu)

🐦 [@jlacomis](https://twitter.com/jlacomis)

🌐 [jeremylacomis.com](http://jeremylacomis.com)