
DEWOLF: IMPROVING DECOMPILATION BY LEVERAGING USER SURVEYS

Steffen Enders, Eva-Maria C. Behner, Niklas Bergmann, Mariia Rybalka,
Elmar Padilla, Er Xue Hui, Henry Low, Nicholas Sim

steffen.enders@fkie.fraunhofer.de

2023-03-03 | BAR Workshop, San Diego



CodeBrowser: ghidra_random_files:/ed01ebfbc9eb5bbea545af4d01bf5f1071661840480439c6e5babe8e080e41aa/ed01ebfbc9eb5bbea545af4d01bf5f1071661840480439c6e5babe8e080e41aa

File Edit Analysis Graph Navigation Search Select Tools Window Help

Program Trees

- ed01ebfbc9eb5bbea545af4d01bf5f1071661840480439c6e5babe8e080e41aa
 - Headers
 - .text
 - .rdata
 - .data

Symbol Tree

- > f __allmul
- > f entry
- > FUN_0040
- > Unwind@004079
- > Labels
- > Classes
- > Namespaces

Data Type Manager

Data Types

- > BuiltInTypes
- > ed01ebfbc9eb5bbea545af4d01bf5f1071661840480439c6e5babe8e080e41aa
- > windows_vs12_32

Decompile: FUN_004014a6 - (ed01ebfbc9eb5bbea545af4d01bf5f1071661840480439c6e5babe8e080e41aa)

```

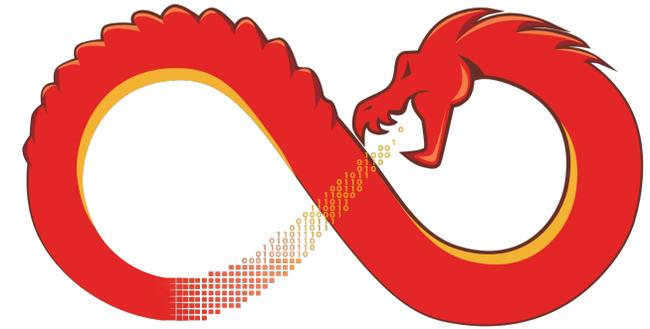
36  uStack_23b = 0;
37  uStack_239 = 0;
38  local_244 = 0;
39  local_20[0] = 0;
40  local_8 = 0;
41  ExceptionList = &local_14;
42  hFile = CreateFileA(param_1,0x80000000,1,(LPSECURITY_ATTRIBUTES)0x0,3,0,(HANDLE)0x0);
43  if (hFile != (HANDLE)0xffffffff) {
44    GetFileSizeEx(hFile,(PLARGE_INTEGER)&local_28);
45    if ((local_24 < 1) && ((local_24 < 0 || (local_28 < 0x6400001)))) {
46      iVar2 = (*DAT_0040f880)(hFile,&local_240,8,local_20,0);
47      if (iVar2 != 0) {
48        iVar2 = memcmp(&local_240,s_WANACRY!_0040eb7c,8);
49        if (iVar2 == 0) {
50          iVar2 = (*DAT_0040f880)(hFile,&local_248,4,local_20,0);
51          if ((iVar2 != 0) && (local_248 == 0x100)) {
52            iVar2 = (*DAT_0040f880)(hFile,*(&local_248 + 0x4c8),0x100,local_20,0);
53            if (iVar2 != 0) {
54              iVar2 = (*DAT_0040f880)(hFile,&local_244,4,local_20,0);
55              if (iVar2 != 0) {
56                iVar2 = (*DAT_0040f880)(hFile,&local_238,8,local_20,0);
57                if (((iVar2 != 0) && ((int)local_234 < 1)) &&
58                    (((int)local_234 < 0 || (local_238 < 0x6400001)))) {
59                  iVar2 = FUN_004019e1((void *)(&local_234 + 4),*(void **)((int)this + 0x4c8),
60                                local_248,local_230,&local_30);
61                }
62                if (iVar2 != 0) {
63                  FUN_00402a76((void *)(&local_234 + 0x54),local_230,(uint *)PTR_DAT_0040f578,
64                            local_30,(byte *)0x10);
65                  local_2c = (byte *)GlobalAlloc(0,local_238);
66                  if (local_2c != (byte *)0x0) {
67                    iVar2 = (*DAT_0040f880)(hFile,*(&local_2c + 0x4c8),local_28,
68                                      local_20,0);
69                    pbVar1 = local_2c;
70                    if (((iVar2 != 0) && (local_20[0] != 0)) &&

```

004014ef FUN_004014a6 STOSW ES:EDI

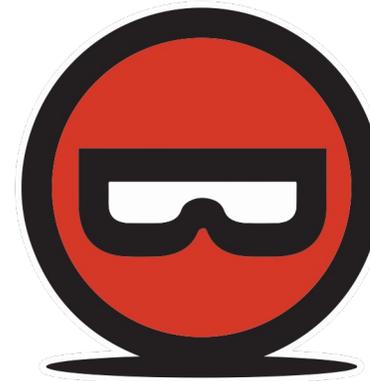
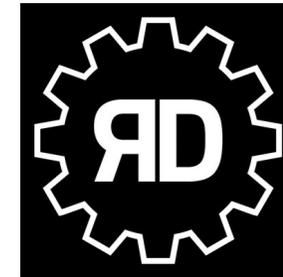
Related Work / Decompilers

- Essential for many binary analysis approaches
- Each decompiler has a different focus!



GHIDRA

USECAP/dream



- There is still plenty of room for improvements!

<https://github.com/revng>

<https://www.qbssoftware.com/jeb-decompiler.html>

<https://mobile.twitter.com/retdec>

<https://ghidra-sre.org>

<https://hex-rays.com>

<https://github.com/USECAP/dream>

<https://binary.ninja>

Decompilation can be done differently!

Output 1

```
undefined4 dividability_rules(void) {
    int iVar1; int local_10 [3];
    printf("Enter a number: ");
    __isoc99_scanf(&DAT_0804c025,local_10);
    iVar1 = local_10[0];
    printf("A number is dividable by %d:\n",local_10[0]);
    if (local_10[0] == 0x7d) {
        printf("last three digits dividable by 125",iVar1);
        return 0; }
    if (local_10[0] < 0x7e) {
        if (local_10[0] < 0xb) {
            switch(local_10[0]) {
                case 0:
                    printf("not possible",iVar1);
                    return 0;
                    ...
                case 10:
                    printf("last digit is 0",iVar1);
                    return 0;
            }
        } else {
            if (local_10[0] == 100) {
                printf("the last two digits are 0",iVar1);
                return 0;
            }
        }
    }
    printf("we have no rule",iVar1); return 0; }
```

Output 2

```
int dividability_rules() {
    int var_0; int * var_1;
    printf("Enter a number: ");
    var_1 = &var_0;
    __isoc99_scanf("%d", var_1);
    printf("A number is dividable by %d:\n", var_0);
    switch(var_0) {
        case 0x0:
            printf("not possible");
            break;
            ...
        case 0xa:
            printf("last digit is 0");
            break;
        case 100:
            printf("the last two digits are 0");
            break;
        case 125:
            printf("last three digits dividable by 125");
            break;
        default:
            printf("we have no rule");
    }
    return 0;
}
```

User Surveys

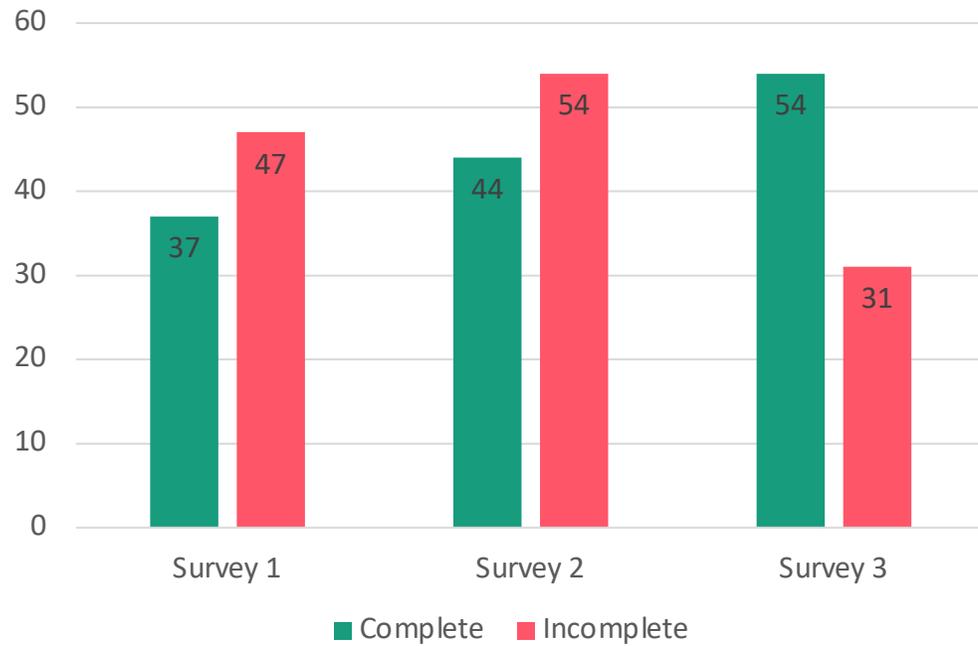
Decompiler Survey Load unfinished survey Exit and clear survey

* Please consider the following decompiled function and answer the questions below. Feel free to use the editor as you would normally do (e.g. by renaming variables).

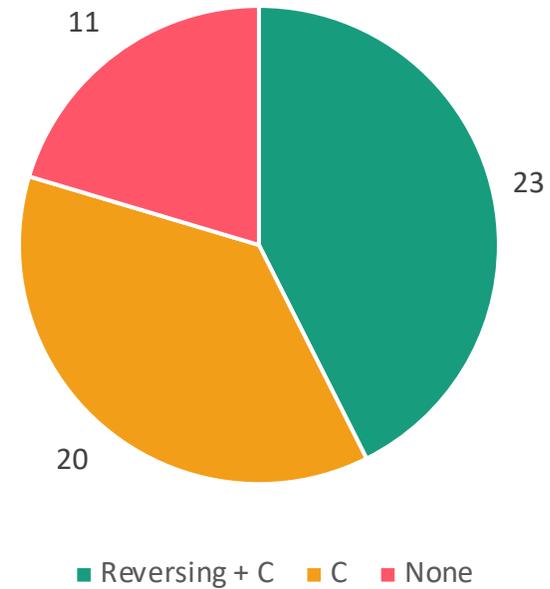
```
1 unsigned long A(int arg1, int arg2) {
2     unsigned long var_0;
3
4     if (arg1 == 0) {
5         var_0 = arg2;
6     }
7     else {
8         var_0 = arg2;
9         while(true) {
10            arg2 = (arg1 + -0x1) & 0xffffffff;
11            if ((unsigned int)var_0 == 0) {
12                if (arg2 == 0) {
13                    var_0 = 0x1;
14                    break;
15                }
16                arg1 = arg2;
17                var_0 = 0x1;
18                continue;
19            }
20            var_0 = A(arg1, ((unsigned int) var_0) + 0xffffffff);
21            if (arg2 == 0) {
22                break;
23            }
24            arg1 = arg2;
25        }
26    }
```

User Surveys Metadata

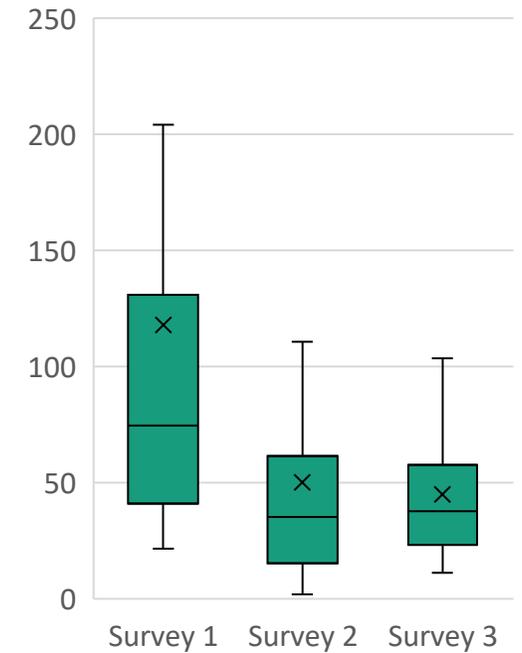
Survey Respondents



Respondents Skills (Survey 3)



Time (min)



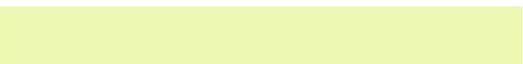
User Survey 1

- Research Questions:
 - What decompilation aspects should be improved to enhance manual analysis?
 - What readability aspects are important to users?
 - What are the limitations of current approaches and the state-of-the-art?

Example 1

```
int test1() {
    int var_0;

    var_0 = rand();
    rand();
    switch(var_0) {
    case 1:
        var_0 = 0;
        break;
    case 5:
        var_0 = var_0 + 1;
    case 10:
        var_0 = var_0 << 1;
        break;
    }
    return var_0;
}
```



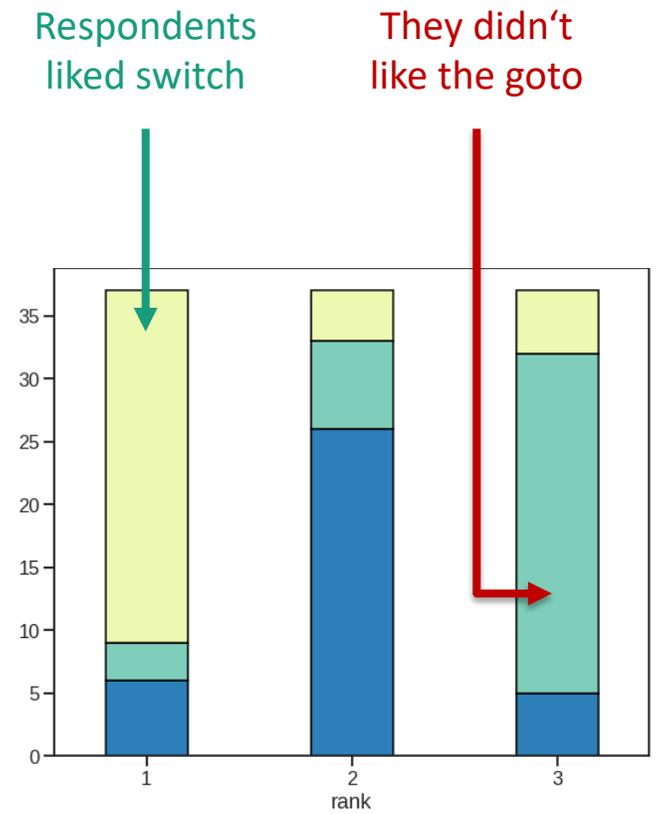
```
int test1(void) {
    int var_0;

    var_0 = rand();
    rand();
    if (var_0 != 10) {
        if (10 < var_0) {
            return var_0;
        }
        if (var_0 == 1) {
            return 0;
        }
        if (var_0 != 5) {
            return var_0;
        }
        var_0 = 6;
    }
    return var_0 << 1;
}
```



```
int test1() {
    int var_0;

    var_0 = rand();
    rand();
    if ( var_0 == 10 ) {
        goto Label_1;
    }
    if ( var_0 <= 10 ) {
        if ( var_0 == 1 ) {
            return 0;
        }
        if ( var_0 == 5 ) {
            var_0 = 6;
        }
    }
Label_1:
    var_0 *= 2;
    return var_0;
}
```



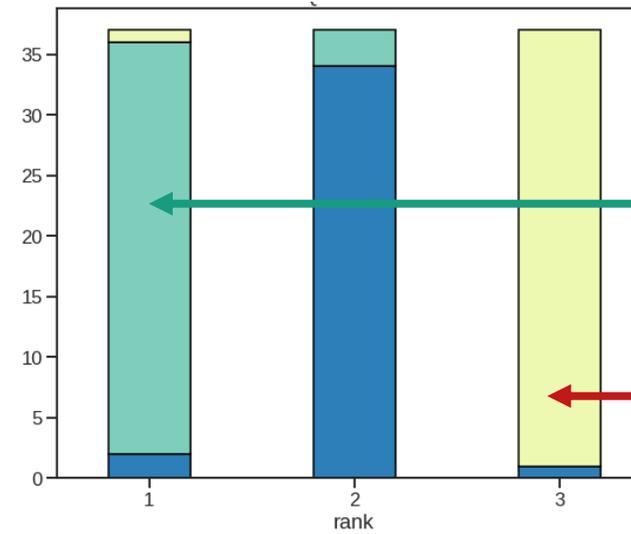
Example 2

```
int test3() {
    int exit_4;
    int var_0;
    int var_1;
    int * var_2;
    int var_3;
    int var_4;

    printf("Please enter a number:");
    var_2 = &(var_0);
    __isoc99_scanf(0x804a01f, var_2);
    while(var_0 <= 99) {
        var_4 = rand();
        var_3 = 0;
        while(true) {
            if (var_3 >= var_4) {
                exit_4 = 0;
                break;
            }
            var_1 = var_3 + var_0;
            if (var_3 + var_0 > 100) {
                var_0 = var_1;
                exit_4 = 1;
                break;
            }
            var_3 = var_3 + 1;
            var_0 = var_1;
        }
        if (exit_4 != 0) {
            break;
        }
    }
    var_3 = printf("The result is %d\\n", var_0);
    return var_3;
}
```

```
void test3(void) {
    int var_0;
    int var_2;
    int var_1;

    printf("Please enter a number:");
    __isoc99_scanf(0x0804a01f,&var_0);
    while (var_0 < 100) {
        var_2 = rand();
        var_1 = 0;
        while (var_1 < var_2) {
            var_0 = var_1 + var_0;
            if (100 < var_0) {
                goto Label_1;
            }
            var_1 = var_1 + 1;
        }
    }
Label_1:
    printf("The result is %d\\n",var_0);
    return;
}
```



+ Short & Clean
+ For-Loop

- Long & Complex
- Too many Variables

```
int test3() {
    int var_0;
    int var_1;
    int i;

    printf("Please enter a number:");
    __isoc99_scanf("%d", &var_0);
    while ( var_0 <= 99 ) {
        var_1 = rand();
        for ( i = 0; i < var_1; ++i ) {
            var_0 += i;
            if ( var_0 > 100 ) {
                return printf("The result is %d\\n", var_0);
            }
        }
    }
    return printf("The result is %d\\n", var_0);
}
```

Survey 1 TL;DR

- Readability depends on user preference *and* the given analysis task
 - Users crave configurable decompilers
- Readability > Assembly Structure
- Identified many aspects with room for improvements
 - Instruction Idioms, Switch Reconstruction, ...

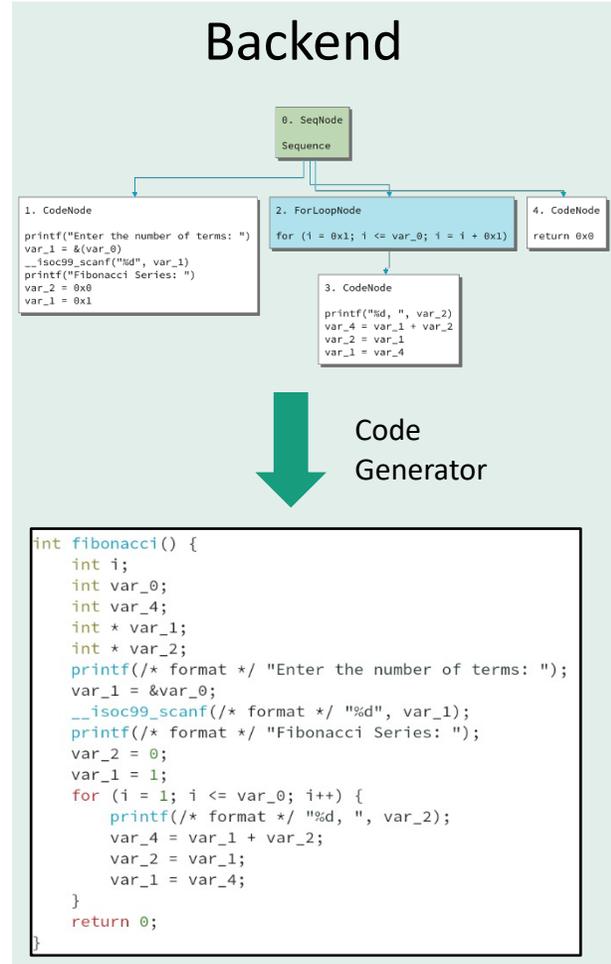
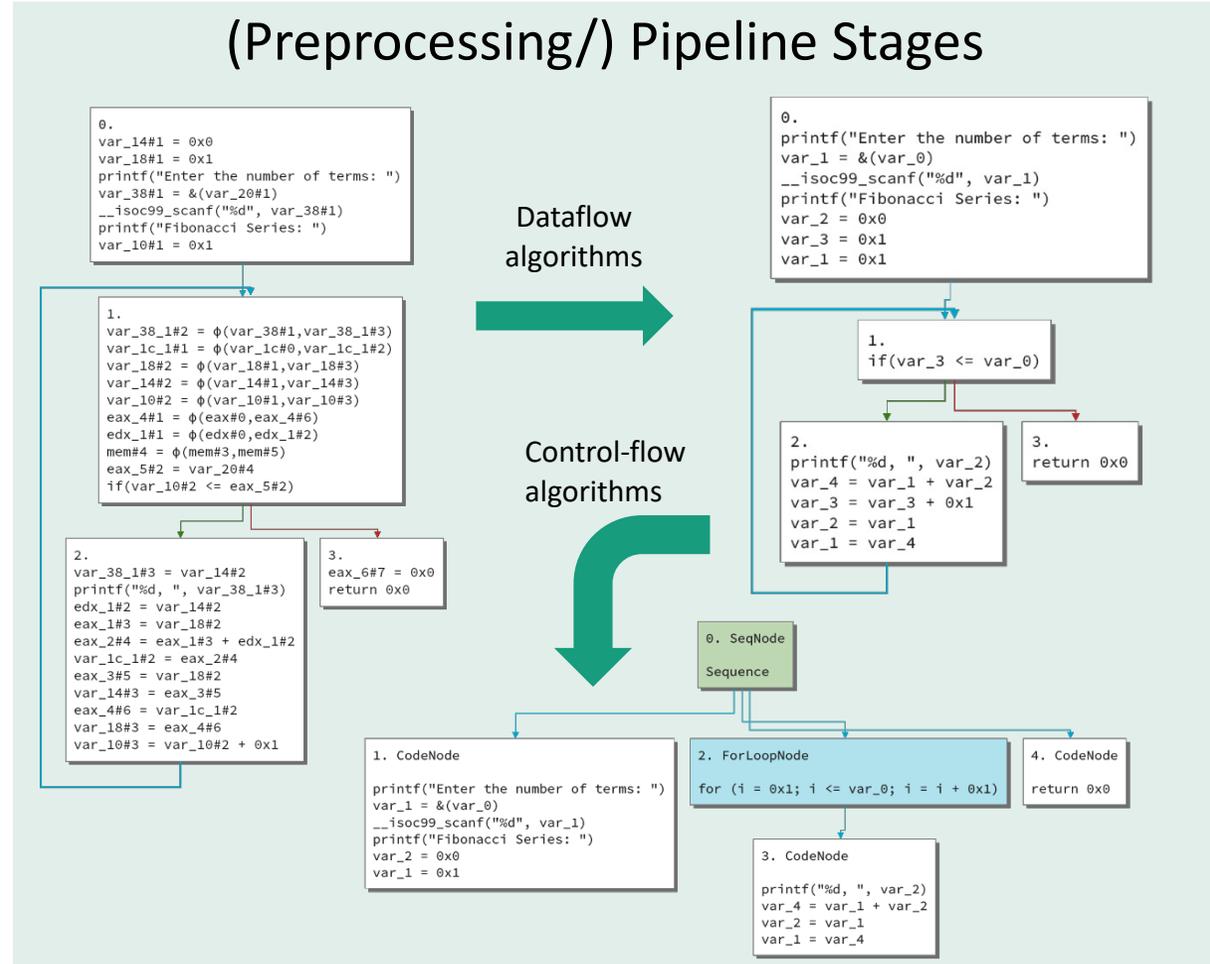
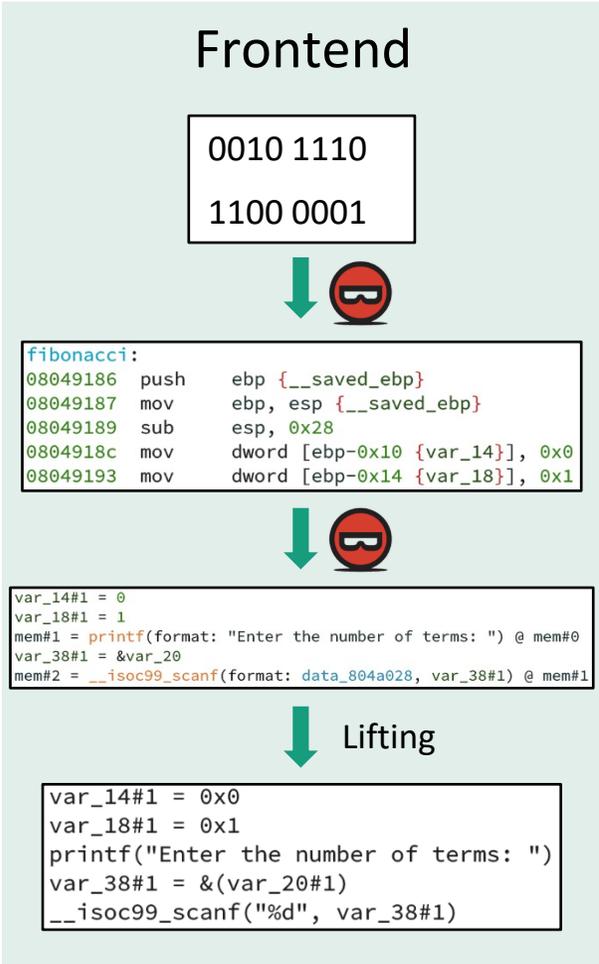
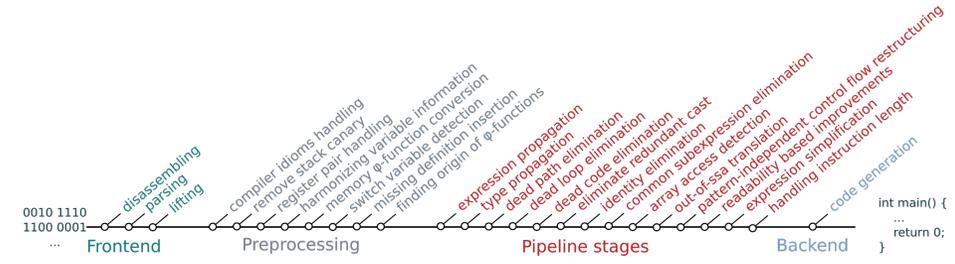
Contributions

- Various improvements for DREAM to get academia “back on track” with commercial state-of-the-art decompilers
- **dewolf Decompiler**
 - highly configurable, easily expandable, and open-source
 - Available in Binary Ninja’s Plugin Manager and on dogbolt.org
- We publish all survey results [1]



[1] <https://github.com/steffenenders/dewolf-surveys>

dewolf Overview



dewolf Improvements over DREAM

Switch Variable Detection
For-loop Recovery
Elimination of Dead Paths and Loops
Custom Logic Engine



Subexpression Elimination
Instruction Idiom Handling
Elimination of Redundant Casts
Improved Out-of-SSA



Constant Representation
Array Access Detection
Continue in Loops



```
int64_t xor(int64_t arg1, int64_t arg2, int64_t arg3, int64_t arg4) {
    printf("%lu", arg2);
    int32_t var_c = 0x0L;
    while(true) {
        int64_t rax_14 = var_c /* sx None */;
        if ((arg2 <= var_c /* sx None */) {
            break;
        }
        char* rax_4 = (arg1 + var_c /* sx None */);
        uint64_t rsi_1 = *rax_4 /* zx char */ /* zx uint32_t */;
        int64_t rax_6 = var_c /* sx None */;
        int64_t rdx_1 = 0x0L;
        char* rax_9 = (arg3 + ((rdx_1:rax_6) % arg4));
        uint64_t rcx = *rax_9 /* zx char */ /* zx uint32_t */;
        int64_t rax_12 = (arg1 + var_c /* sx None */);
        uint64_t rdx_4 = (rsi_1.esi ^ rcx.ecx) /* zx None */;
        *rax_12 = rdx_4.dl;
        int32_t var_c = (var_c + 0x1L);
    }
    return rax_14;
}
```

vs.

```
long xor(void * arg1, long arg2, long arg3, long arg4) {
    int i;
    int var_1;
    void * var_0;
    printf(/* format */ "%lu\n", arg2);
    for (i = 0; arg2 > i; i++) {
        var_0 = arg1 + i;
        *var_0 = *var_0 ^ *(arg3 + (i + (0L << 0x40)) % arg4);
    }
    var_1 = i;
    return var_1;
}
```

User Survey 2+3 Research Questions

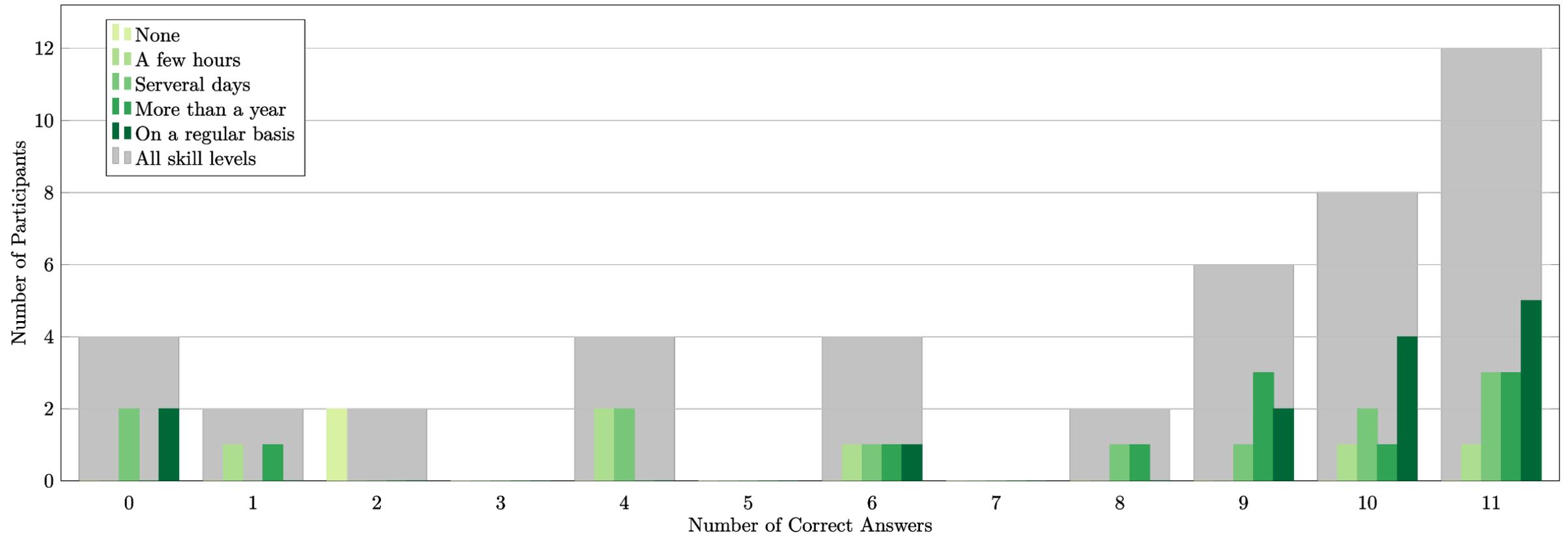
- Research Questions:
 - Does dewolf produce readable & comprehensible output to qualify as a base for future research?
 - (Can dewolf exceed state-of-the-art decompilers in certain cases?)
 - To what extent do respondents tolerate diverging from the assembly during decompilation?
 - What are further areas to improve readability in future research?

Comprehension Questions

- What is the purpose of the function?
- What TLDs does the DGA utilize?
- What is the most used TLD?
- Which of the provided domains could be generated by the DGA?
- How many different domains can be generated by the DGA?

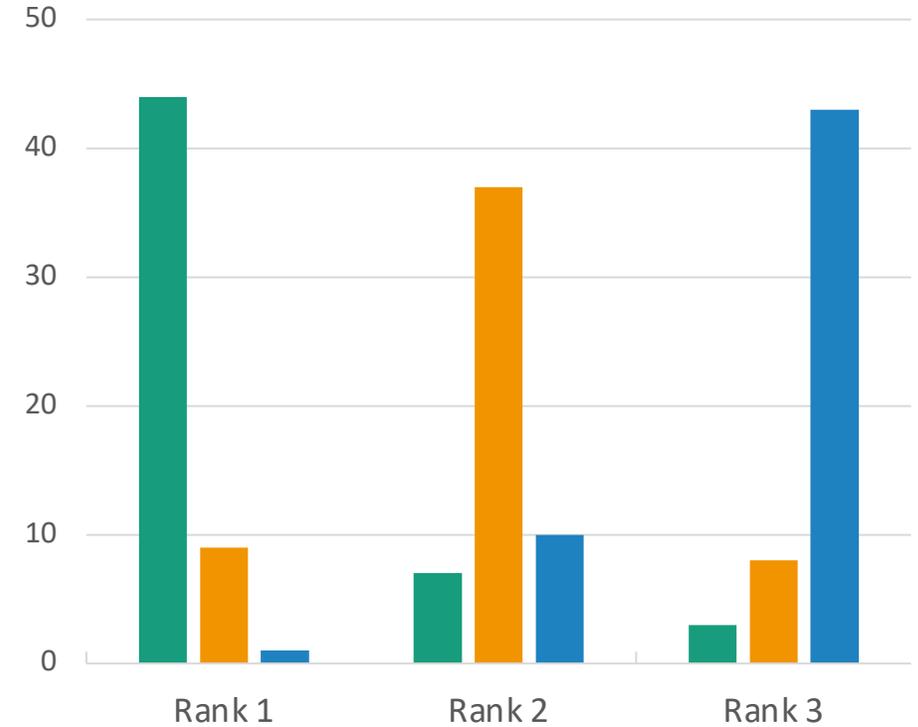
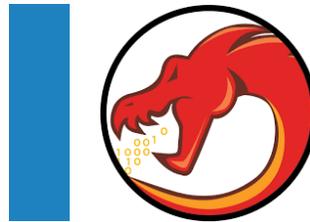
```
int sub_401000() {
    char i;
    char * var_0;
    char var_1;
    char var_3;
    int var_4;
    var_0 = malloc(16);
    var_1 = GetTickCount();
    GetSystemTime(lpSystemTime: var_0);
    for (i = 0; i < 8; i++) {
        var_3 = var_0[i];
        var_0[i] = ((unsigned int)(var_3 ^ var_1) % 24 & 0xff) + 'a';
    }
    var_4 = var_0 + 8;
    *var_4 = 0x6d6f632e;
    *(var_4 + 4) = 0x0;
    switch((((int) var_1) % 0x8) - 0x1) {
    case 0:
    case 5:
        *var_4 = *var_4 ^ 0x6d001700;
        break;
    case 1:
        *var_4 = *var_4 ^ 0x190a0d00;
        break;
    case 4:
    case 6:
        *var_4 = *var_4 ^ 0x6d1a1100;
        break;
    }
    return var_0;
}
```

Comprehension Results Divided by C-Skills (self assessment)



Decompiler Comparison

```
int convert_binary_to_hex()
{
    long long binary;
    char hex[65] = "";
    int rem;
    printf("Enter any binary number: ");
    scanf("%lld", &binary);
    int original_binary = binary;
    while(binary > 0)
    {
        rem = binary % 10000;
        switch(rem)
        {
            case 0:
                strcat(hex, "0");
                break;
            case 1:
                strcat(hex, "1");
                break;
            case 10:
                ...
            case 1111:
                strcat(hex, "F");
                break;
        }
        binary /= 10000;
    }
    printf("Binary number: %lld\n", original_binary);
    printf("Hexadecimal number: %s", hex);
    return 0;
}
```



What participants (dis)liked in state-of-the-art decompilers

Positive

Switch-Recovery
Low Nesting Depth
for-loops, if possible
Good Type Recovery
Less Casts
Array Access Detection
Constant Annotations

vs.

Negative

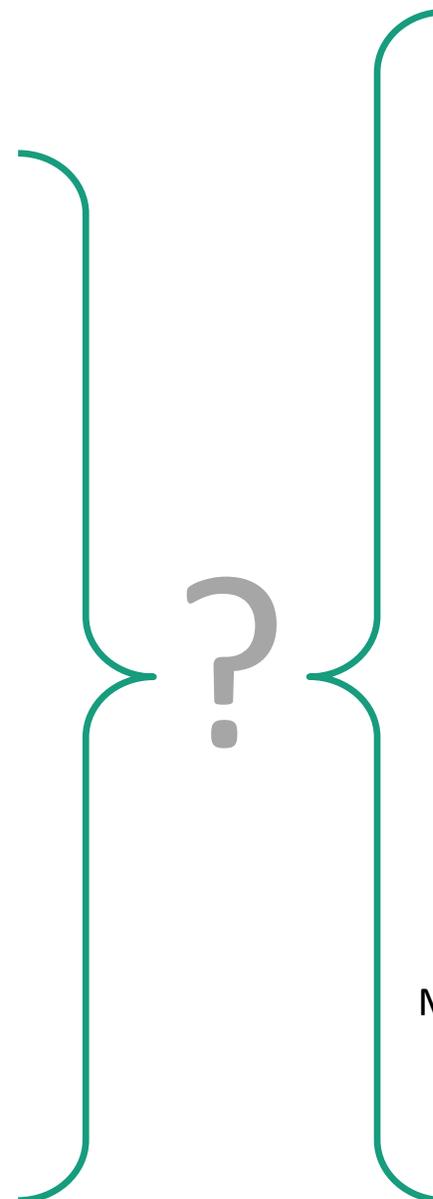
Deeply Nested If/Else
High Nesting Depth
While-loops instead of for-loops
Bad Type Recovery
Explicit Casts
Too many unnecessary variables
String Parameters

Introducing Copies during Restructuring?

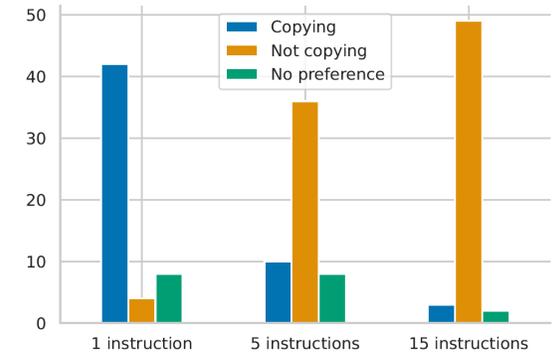
```
return var;
```

```
y = x + 5;
z = bar(y);
printf("Some text to print");
y += z;
return y;
```

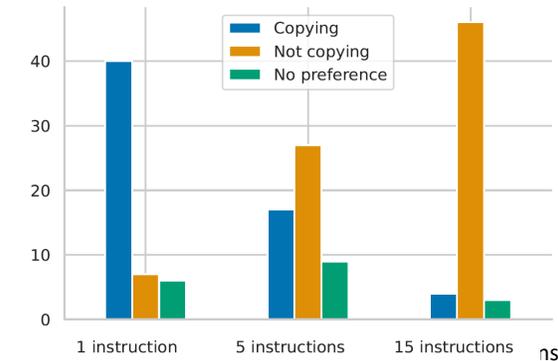
```
y = x + 5;
z = bar(y);
printf("Enter a number larger than %d: \n", z);
scanf("%d", &numb1);
printf("Enter a number smaller than %d: \n", z);
scanf("%d", &numb2);
y += z;
printf("Enter a number larger than %d: \n", y);
scanf("%d", &numb3);
printf("Enter a number smaller than %d: \n", y);
scanf("%d", &numb4);
diff_1 = numb1 - numb2;
diff_2 = numb3 - numb4;
printf("The two differences are %d and %d: \n", diff_1, diff_2);
return diff_1 + diff_2;
```



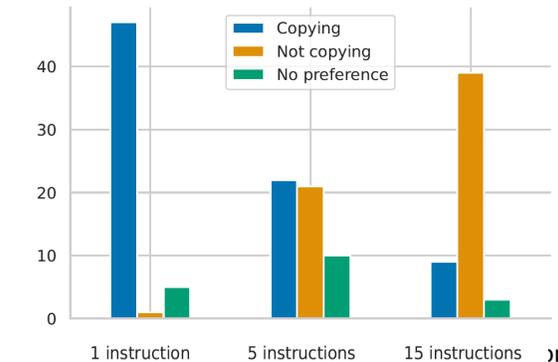
If-Else



Multi-Exit



Multi-Entry



TL;DR

- Future decompilers (focusing on manual analysis) should:
 - Favor readability over sticking to the assembly
 - Be highly configurable!

- If you are into decompilation, you should:
 - Check out dewolf (great base for new research, imho)
 - Read the paper and survey results



Q&A

steffen.enders@fkie.fraunhofer.de



@steffenenders_



<https://github.com/fkie-cad/dewolf>

