# No More Gotos: Decompilation Using Pattern-Independent Control-Flow Structuring and Semantics-Preserving Transformations

Khaled Yakdan[1]    Sebastian Eschweiler[2]    Elmar Gerhards-Padilla[2]
Matthew Smith[1]

[1]University of Bonn, Germany

[2]Fraunhofer FKIE, Germany

NDSS 2015

universität**bonn**

# Agenda

universität**bonn**

```
0101010101010101010100
0101010101010101010100
0101010101010101010100
0101010101010101010100
0101010101010101010100
```

**Binary code**

# Motivation
## Decompilation in Security

**Source code**

```
int f(int a){
    int i = 0;
    for(; i < a ; i++)
        ...
}
```
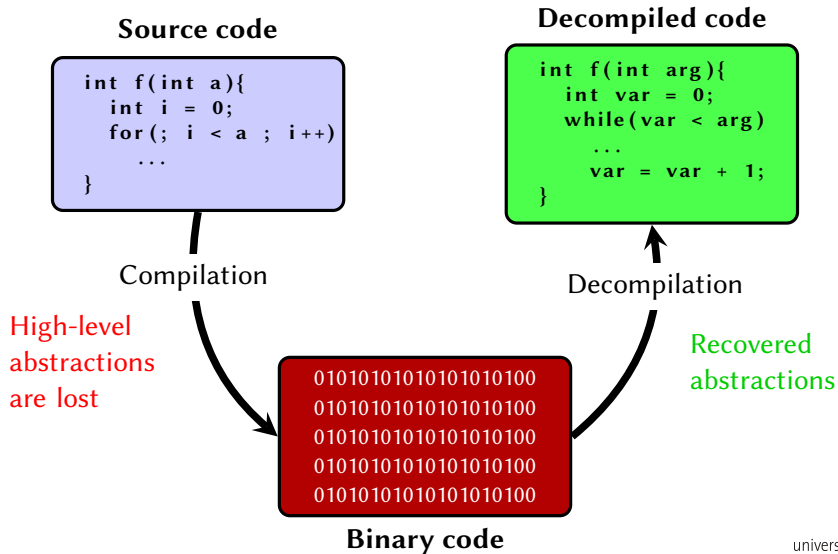
**Decompiled code**

```
int f(int arg){
    int var = 0;
    while(var < arg)
        ...
        var = var + 1;
}
```

Compilation

Decompilation

High-level abstractions are lost

Recovered abstractions

```
01010101010101010100
01010101010101010100
01010101010101010100
01010101010101010100
01010101010101010100
```

**Binary code**

universität**bonn**

- Manual reverse engineering

- Manual reverse engineering
- Apply source-based techniques to binary code

universität**bonn**

- Manual reverse engineering
- Apply source-based techniques to binary code
  - Find vulnerabilities, bugs

# Motivation
## Decompilation in Security

- Manual reverse engineering
- Apply source-based techniques to binary code
  - Find vulnerabilities, bugs
  - Taint tracking

- Manual reverse engineering
- Apply source-based techniques to binary code
  - Find vulnerabilities, bugs
  - Taint tracking

### Goal: Enhanced Structuredness

Effective control flow structure recovery to improve readability and enhance program analysis

# Agenda

universität**bonn**

# Control Flow Structuring

```
int f(int a, int b){
  int sum = 0;
  if(a < b){
    for(int i = a; i < b; i++)
      sum += i;
  }
  return sum;
}
```

Structured code

```
int f(int a, int b){
  int sum = 0;
  if(a >= b)
    goto Label_2;
  int i = a;
Label_1:
  if(i >= b)
    goto Label_2;
  sum += i;
  i++;
  goto Label_1;
Label_2:
  return sum;
}
```

Unstructured code

universität**bonn**

State of the art: **Structural Analysis** [Sharir80]

- Pattern-matching using a predefined set of region schemas (patterns)

State of the art: **Structural Analysis** [Sharir80]

- Pattern-matching using a predefined set of region schemas (patterns)
- Use goto statements if no match is found

State of the art: **Structural Analysis** [Sharir80]

- Pattern-matching using a predefined set of region schemas (patterns)
- Use goto statements if no match is found
- Example: Decompiling a P2P Zeus sample with Hex-Rays
  - 1,571 goto for 49,514 LoC
  - **1 goto for each 32 LoC**

# Prior Work on Control-Flow Structuring

Improving vanilla structural analysis to recover more structure

- SESS Analysis [Engel et al., SCOPES 2011]
- Phoenix Decompiler [Schwartz et al., USENIX Security 2013]

universität**bonn**

# Prior Work on Control-Flow Structuring

Improving vanilla structural analysis to recover more structure

- SESS Analysis [Engel et al., SCOPES 2011]
- Phoenix Decompiler [Schwartz et al., USENIX Security 2013]

New control-flow structuring algorithm

**Pattern-Independent Structuring**

**Semantics-Preserving Transformations**

universität**bonn**

# Agenda

universität**bonn**

**IR Trans-formation**

universität**bonn**
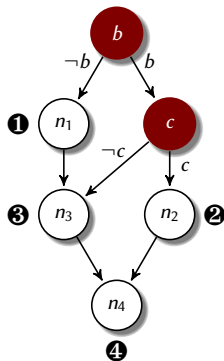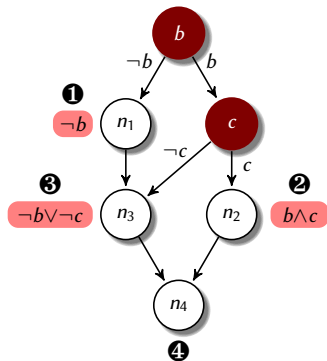
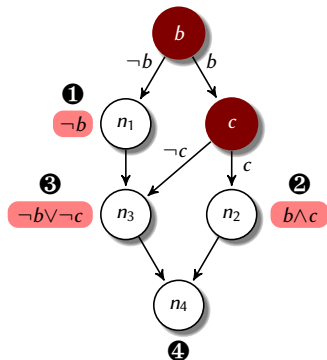# Running Example

- Lexical order

# Acyclic Regions

- Lexical order
- Reaching conditions

# Acyclic Regions

- Lexical order
- Reaching conditions
- Initial AST as a sequence of **if** constructs



```
if (¬b)
    n₁
if (b ∧ c)
    n₂
if (¬b ∨ ¬c)
    n₃
n₄
```
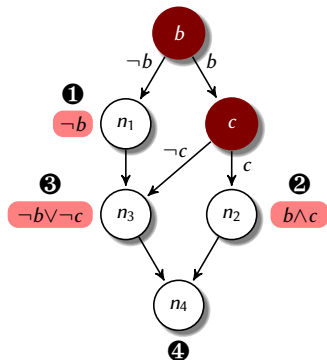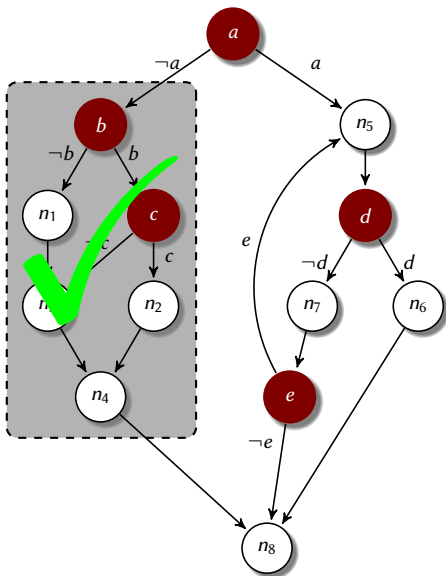
# Acyclic Regions

- Lexical order
- Reaching conditions
- Initial AST as a sequence of **if** constructs
- Refine initial AST to find **switch**, **if-else** constructs



```
if (¬b)
    n₁
if (b ∧ c)
    n₂
if (¬b ∨ ¬c)
    n₃
n₄
```

```
if (¬b)
    n₁
if (b ∧ c)
    n₂
else
    n₃
n₄
```
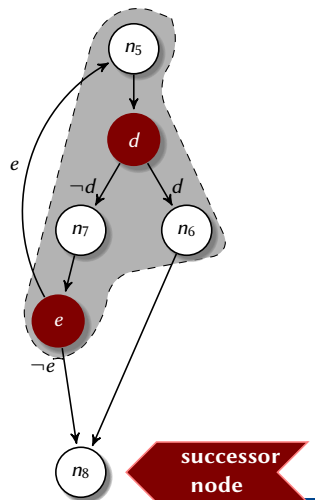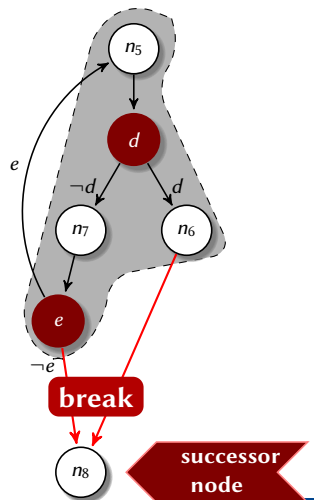
# Cyclic Regions

- Identify loop nodes and successor node

# Cyclic Regions

- Identify loop nodes and successor node
- Replace edges to the successor node by **break** statements

- Identify loop nodes and successor node
- Replace edges to the successor node by **break** statements
- Structure loop body $B_{AST}$

- Identify loop nodes and successor node
- Replace edges to the successor node by **break** statements
- Structure loop body $B_{AST}$
- Initial AST: **while** $(1)$ $\{B_{AST}\}$



```
while (1)
   ...
   if (¬e)
      break
```
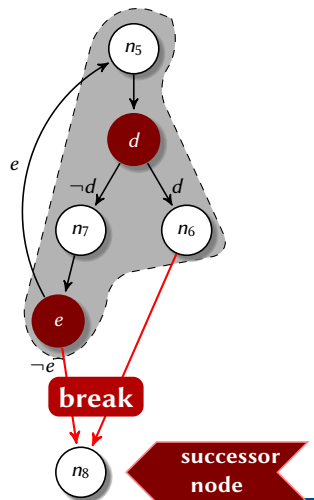
# Cyclic Regions
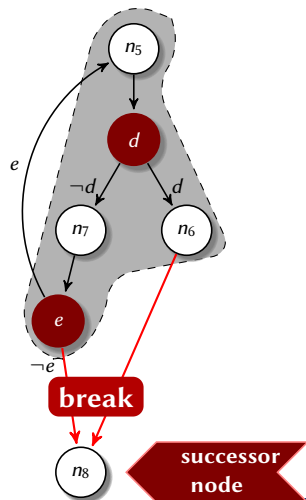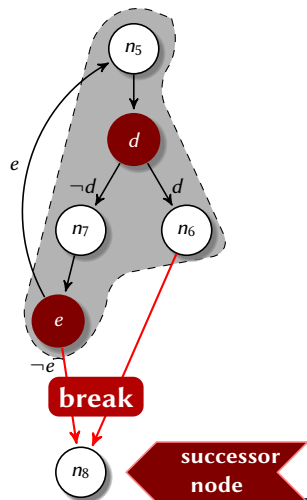
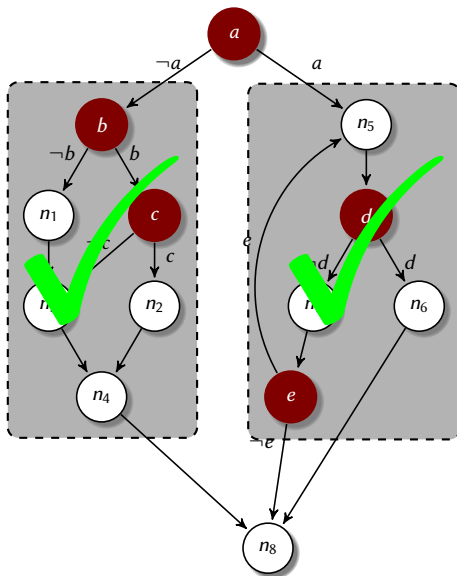- Identify loop nodes and successor node
- Replace edges to the successor node by **break** statements
- Structure loop body $B_{AST}$
- Initial AST: **while** $(1)$ $\{B_{AST}\}$
- Infer loop type and condition

```
while (1)
    . . .
    if (¬e)
        break
```

```
do
    . . .
while (e)
```

# Agenda

01 Motivation

02 Control Flow Structuring

03 The DREAM Decompiler

04 **Results**

05 Conclusion

# Readability Enhancements

```
int f(int a1){
  int v2 = 0;
  while((a1 <= 1 && v2 <= 100)
        || (a1 > 1 && v2 <= 10)){
    printf("inside_loop");
    ++v2;
    --a1;
  }
  printf("loop_terminated");
  return v2;
}
```

Dream

```
signed int f(signed int a1){
  signed int v2;
  v2 = 0;
  while ( a1 > 1 ){
    if ( v2 > 10 )
      goto LABEL_7;
LABEL_6:
    printf("inside_loop");
    ++v2;
    --a1;
  }
  if ( v2 <= 100 )
    goto LABEL_6;
LABEL_7:
  printf("loop_terminated");
  return v2;
}
```
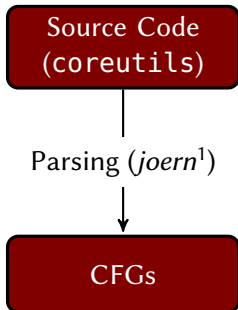
Hex-Rays

universität bonn

Metrics

- Correctness
- Structuredness
- Compactness

universität**bonn**

Source Code
(coreutils)

[1][Yamaguchi et al. IEEE S&P 2014]

universität**bonn**

Source Code
(coreutils)

Parsing (*joern*[1])

CFGs

[1][Yamaguchi et al. IEEE S&P 2014]

universität**bonn**

[1][Yamaguchi et al. IEEE S&P 2014]

Source Code
(coreutils)

Parsing (*joern*[1])

CFGs

Control Flow
Structuring

Structured Code

Passes tests?

---

[1][Yamaguchi et al. IEEE S&P 2014]

| Considered Functions $F$ | $|F|$ | Number of gotos |
|---|---|---|
| Functions after preprocessor | 1,738 | 219 |
| Functions correctly parsed by *joern* [2] | 1,530 | 129 |
| Functions passed tests after structuring | 1,530 | 0 |

[2] Errors have been reported to *joern*'s authors and are fixed in the current release
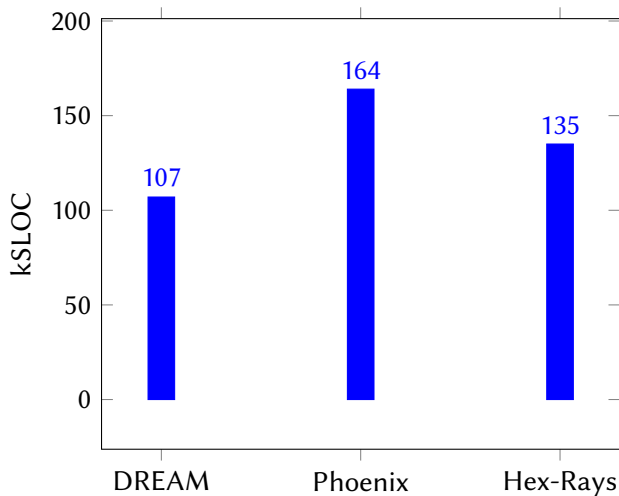
# Structuredness and Compactness

- Tested decompilers
  - DREAM
  - Phoenix (academic state of the art)
  - Hex-Rays (industry state of the art)
- Structuredness
  - Number of gotos
- Compactness
  - Total lines of code
  - Compact functions

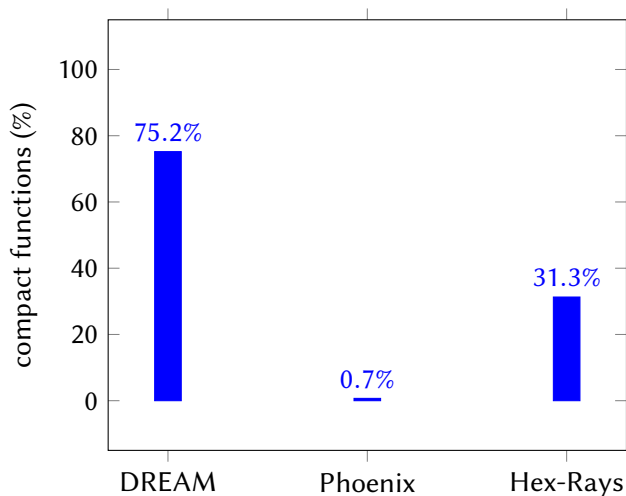universität**bonn**

# Compactness

# Compactness

# Agenda

universität**bonn**

# Conclusion

- Novel control flow structuring algorithm
  - pattern-independent structuring
  - semantics-preserving transformations
- DREAM decompiler
  - goto-free decompiled code
  - compact code
  - good readability

universität**bonn**

# Thank You!

Questions?