

Static Analysis of Virtualization- Obfuscated Binaries

Johannes Kinder

*School of Computer and Communication Sciences
École Polytechnique Fédérale de Lausanne (EPFL), Switzerland*

Virtualization Obfuscation

- Obfuscation
 - *Make code hard to understand for humans and tools*
 - *Popular for protecting benign and malicious code*
- Virtualization Obfuscation
 - *Hide code inside self-contained VM*
 - *Considered one of the strongest obfuscation schemes*
 - *Commercial tools (CodeVirtualizer, VMProtect) and research obfuscator*

Static Analysis vs. Virtualization

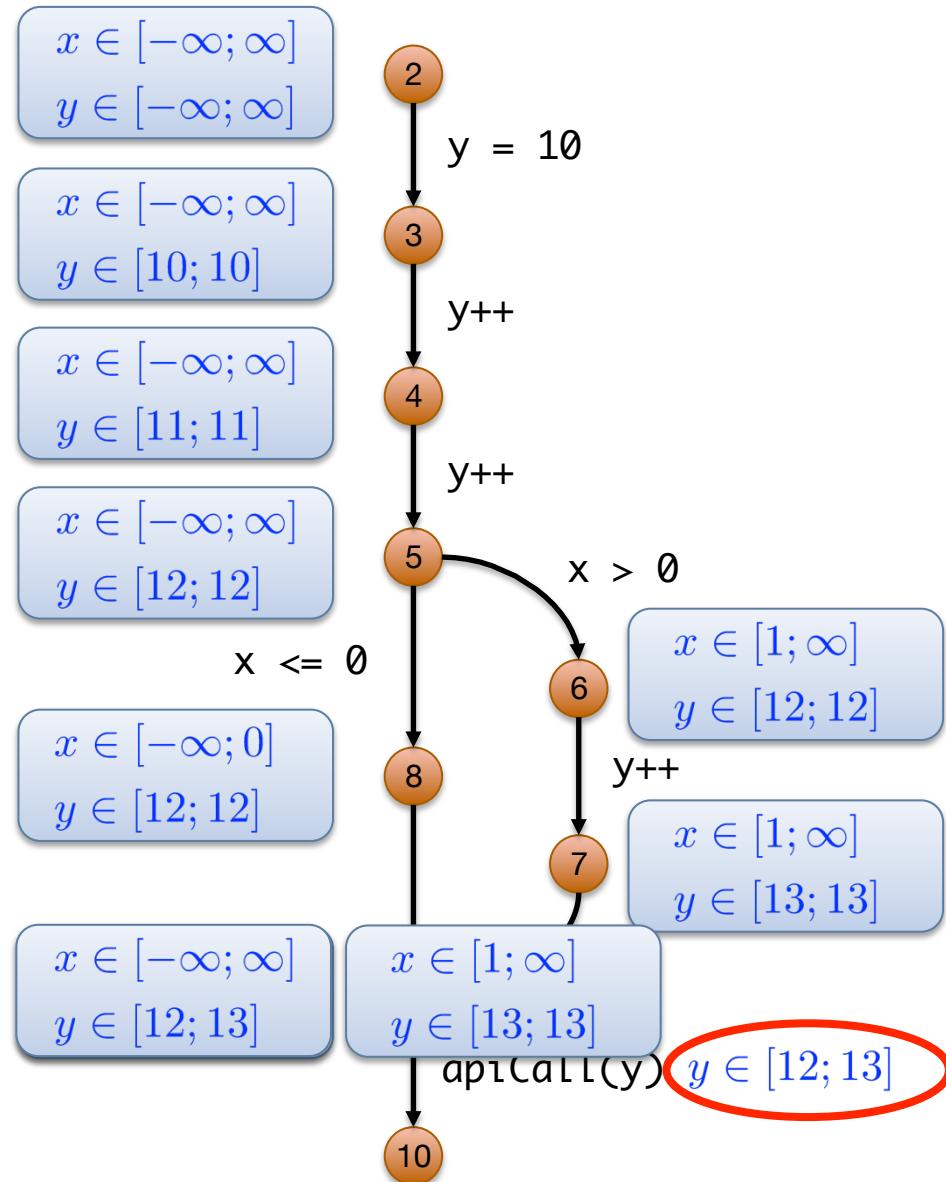
- Static Analysis
 - *Computes program invariants*
 - *Safety proofs, bug finding, malware detection*

What makes static analysis of virtualization-obfuscated code fundamentally hard?

*Can we still make it work and
use it for deobfuscation?*

Static Analysis

```
1 void foo (int x) {  
2     int y = 10;  
3     y++;  
4     y++;  
5     if (x > 0) {  
6         y++;  
7     }  
8     else {}  
9     apiCall(y);  
10 }
```



Virtualization Obfuscation

```
1 void foo (int x) {  
2     int y = 10;  
3     y++;  
4     y++;  
5     if (x > 0) {  
6         y++;  
7     }  
8     else {}  
9     apiCall(y);  
10 }
```

code = { 52, 01, 02, 03, 01, 03, 01, 08,
00, 03, 03, 01, 18, 01, 00 }

data = { 00, 00, 10, 05 }

x
y
conditional
jump distance

Compile to
Bytecode

```
→ 5 int vpc = 0, op1, op2;
  6 while (true) {
  7     switch(code[vpc]) {
  8         case 03: // increment
  9             op1 = code[vpc + 1];
 10            data[op1]++;
 11            vpc += 2;
 12            break;
 13        case 08: // conditional jump
 14            op1 = code[vpc + 1];
 15            op2 = code[vpc + 2];
 16            if (data[op1] <= 0)
 17                vpc += data[op2]
 18            else
 19                vpc += 3;
 20            break;
 21        case 18: // call function
 22            op1 = code[vpc + 1];
 23            apiCall(data[op1]);
 24            vpc += 2;
 25            break;
 26        case 52: // assignment
 27            op1 = code[vpc + 1];
 28            op2 = code[vpc + 2];
 29            data[op1] = data[op2];
 30            vpc += 3;
 31            break;
 32        default: // halt
 33            return;
 34    } // end switch
 35 } // end while
```

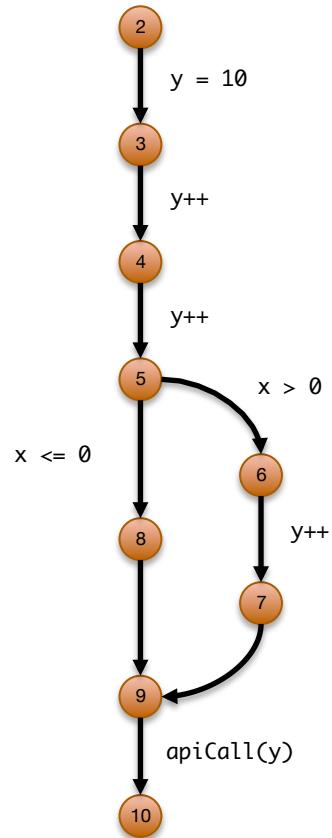
↓

code = { 52, 01, 02, 03, 01, 03, 01, 08,
 00, 03, 03, 01, 18, 01, 00 }

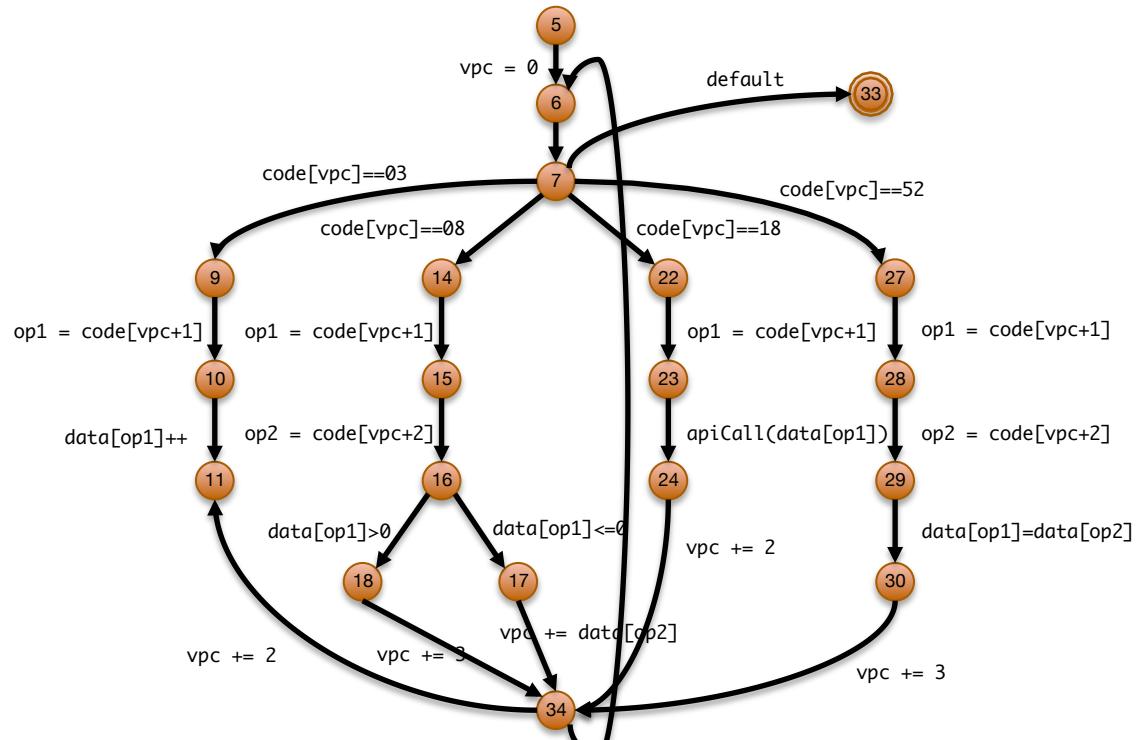
data = { 00, 00, 10, 05 }

x y conditional
 jump distance

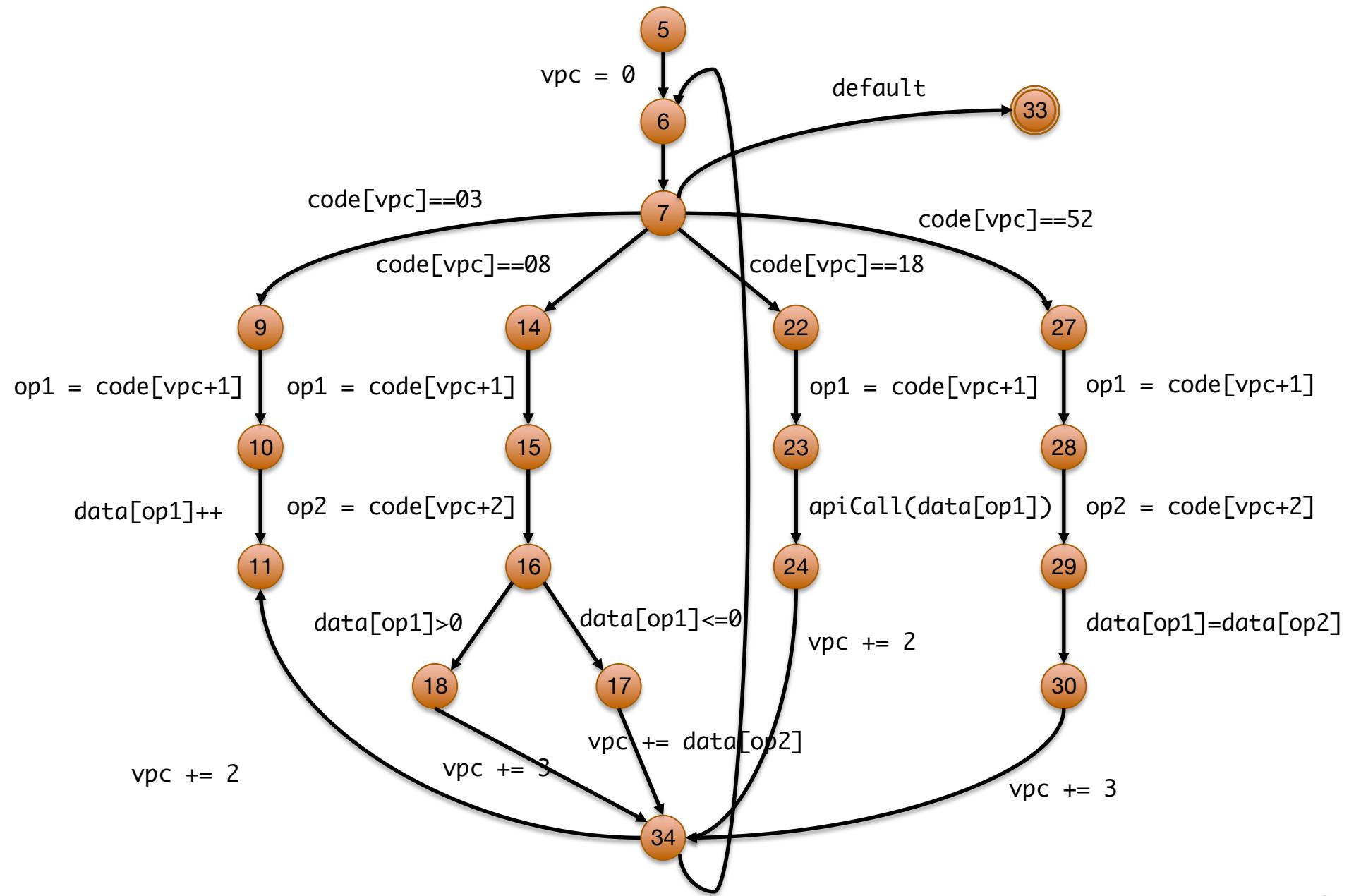
Control Flow Graphs



Original CFG



Obfuscated CFG



```
code = { 52, 01, 02, 03, 01, 03, 01, 08, 00, 03, 03, 01, 18, 01, 00 }
```

```
data = { 00, 00, 10, 05 }
```

d_x d_y d_c d_j

$vpc \in [-\infty; \infty]$
 $d_y \in [0; 0]$

$vpc = 0$

5

6

$vpc \in [0; 7]$ $d_c \in [10; \infty]$
 $d_y \in [0; \infty]$ $d_j \in [5; \infty]$

$code[vpc] == 03$

9

$vpc \in [3; 5]$ $d_c \in [10; 11]$
 $d_y \in [0; 12]$ $d_j \in [5; 6]$

$op1 = code[vpc+1]$

10

$op1 \in [1; 3]$

$data[op1]++$

11

$vpc \in [3; 5]$ $d_c \in [10; 12]$
 $d_y \in [0; 13]$ $d_j \in [5; 7]$

$vpc += 2$

34

$vpc \in [5; 7]$ $d_c \in [10; 12]$
 $d_y \in [0; 13]$ $d_j \in [5; 7]$

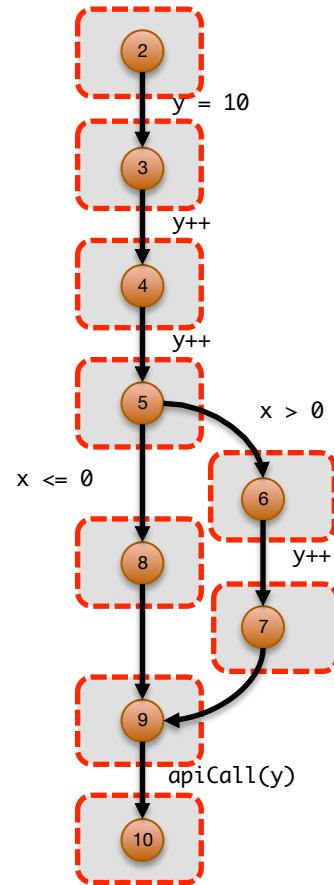
Upper bounds grow to infinity

Weak update

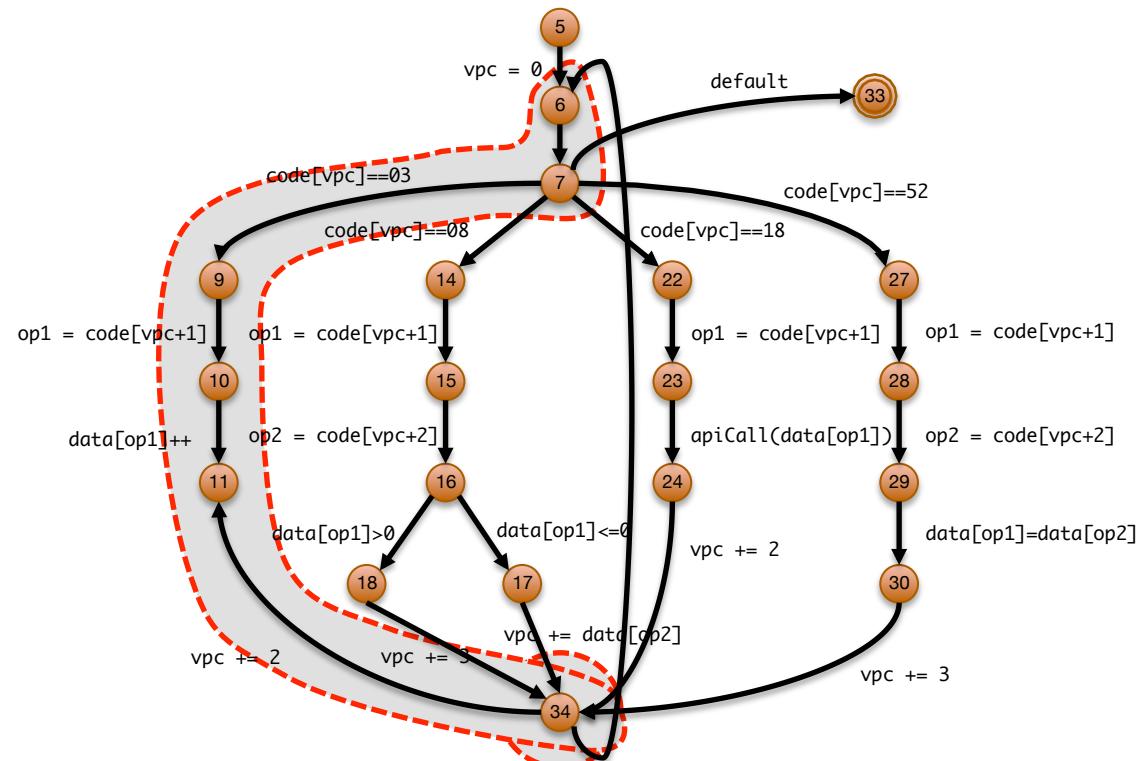
Constant imprecise

Jump distance imprecise

- 1 interpreter case = many original locations
- Interpreter loop head shared among all

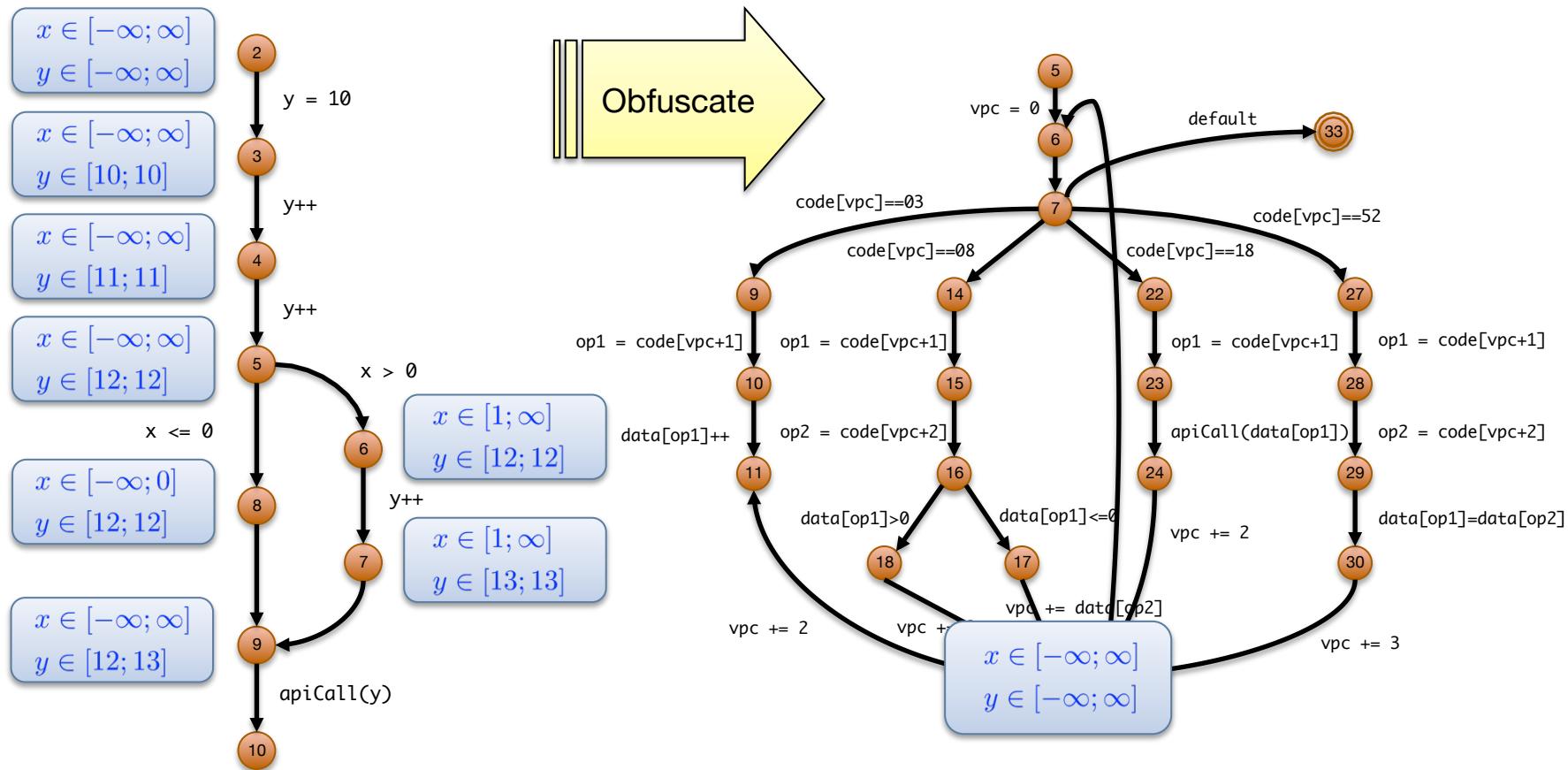


Original CFG



Obfuscated CFG

Domain Flattening



Location Sensitive Analysis → Location Insensitive Analysis

VPC Lifting

- Virtualization flattens one dimension of location
- Idea: track VPC and use as additional dimension
 - *Separate states with differing VPC values*

```
int y = 10;      y++;
```

$$vpc \in [0; 0]$$

$$d_y \in [0; 0]$$

$$vpc \in [3; 3]$$

$$d_y \in [10; 10]$$

- *Join states with equal VPC values*

```
if (...) {}    else {}
```

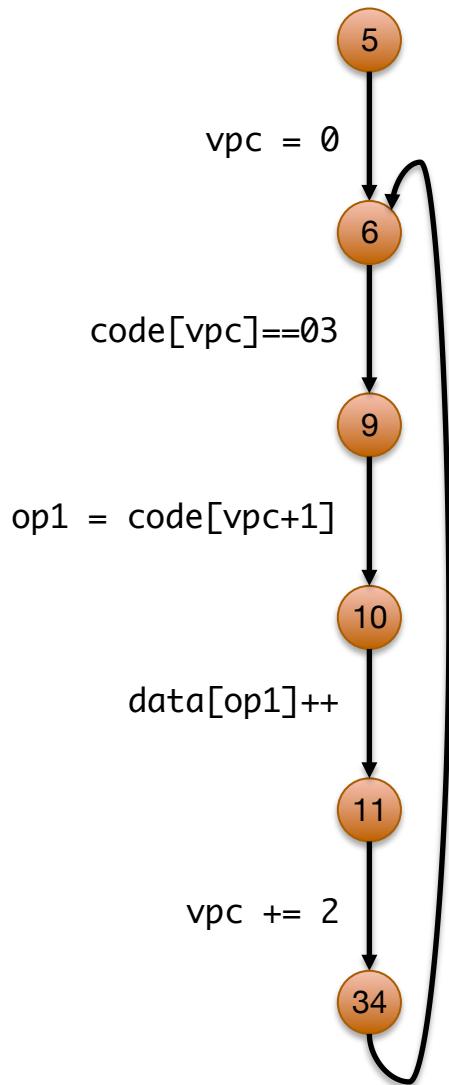
$$vpc \in [12; 12]$$

$$d_y \in [12; 13]$$

$$vpc \in [12; 12]$$

$$d_y \in [13; 13]$$

$y++$	$y++$	$\{ \ y++ \ }$	$\text{apiCall}(y)$
$vpc \in [3; 3]$	$[5; 5]$	$[10; 10]$	$[12; 12]$



$d_x \in [-\infty; \infty]$ $d_y \in [10; 10]$	$d_x \in [-\infty; \infty]$ $d_y \in [11; 11]$	$d_x \in [1; \infty]$ $d_y \in [12; 12]$	$d_x \in [-\infty; \infty]$ $d_y \in [12; 13]$
$d_x \in [-\infty; \infty]$ $d_y \in [10; 10]$	$d_x \in [-\infty; \infty]$ $d_y \in [11; 11]$	$d_x \in [1; \infty]$ $d_y \in [12; 12]$	$d_x \in [-\infty; \infty]$ $d_y \in [12; 13]$
$d_x \in [-\infty; \infty]$ $d_y \in [10; 10]$	$d_x \in [-\infty; \infty]$ $d_y \in [11; 11]$	$d_x \in [1; \infty]$ $d_y \in [12; 12]$	$d_x \in [-\infty; \infty]$ $d_y \in [12; 13]$
$d_x \in [-\infty; \infty]$ $d_y \in [11; 11]$	$d_x \in [-\infty; \infty]$ $d_y \in [12; 12]$	$d_x \in [1; \infty]$ $d_y \in [13; 13]$	$d_x \in [-\infty; \infty]$ $d_y \in [12; 13]$
$d_x \in [-\infty; \infty]$ $d_y \in [11; 11]$	$d_x \in [-\infty; \infty]$ $d_y \in [12; 12]$	$d_x \in [1; \infty]$ $d_y \in [13; 13]$	$d_x \in [-\infty; \infty]$ $d_y \in [12; 13]$

VPC Values as Locations

- Location-sensitive analysis over domain A has domain

$$L \rightarrow A$$

- VPC-sensitive analysis over domain A , for VPC domain V , has domain

$$\underbrace{L \rightarrow V}_{\text{VPC Location}} \rightarrow A$$

“VPC Location”

Reconstructing CFGs

- A program CFG is the set of all feasible transitions between *locations*

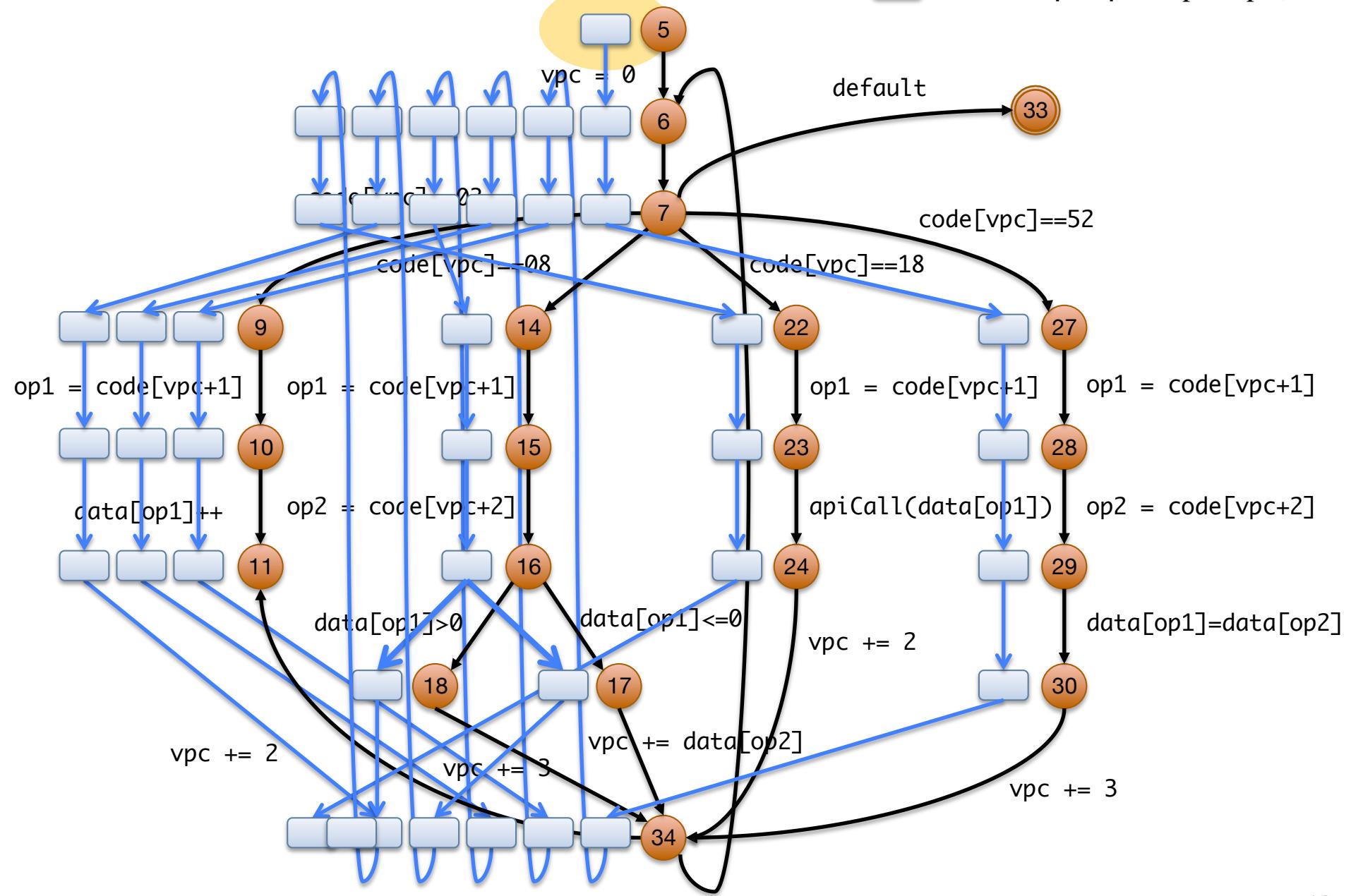
$$CFG = \{(\ell, \ell') \mid (\ell, a) \mapsto (\ell', a'), \ell, \ell' \in L, a, a' \in A\}$$

↑
Abstract transition relation

- A VPC-CFG is the set of all feasible transitions between *VPC locations*

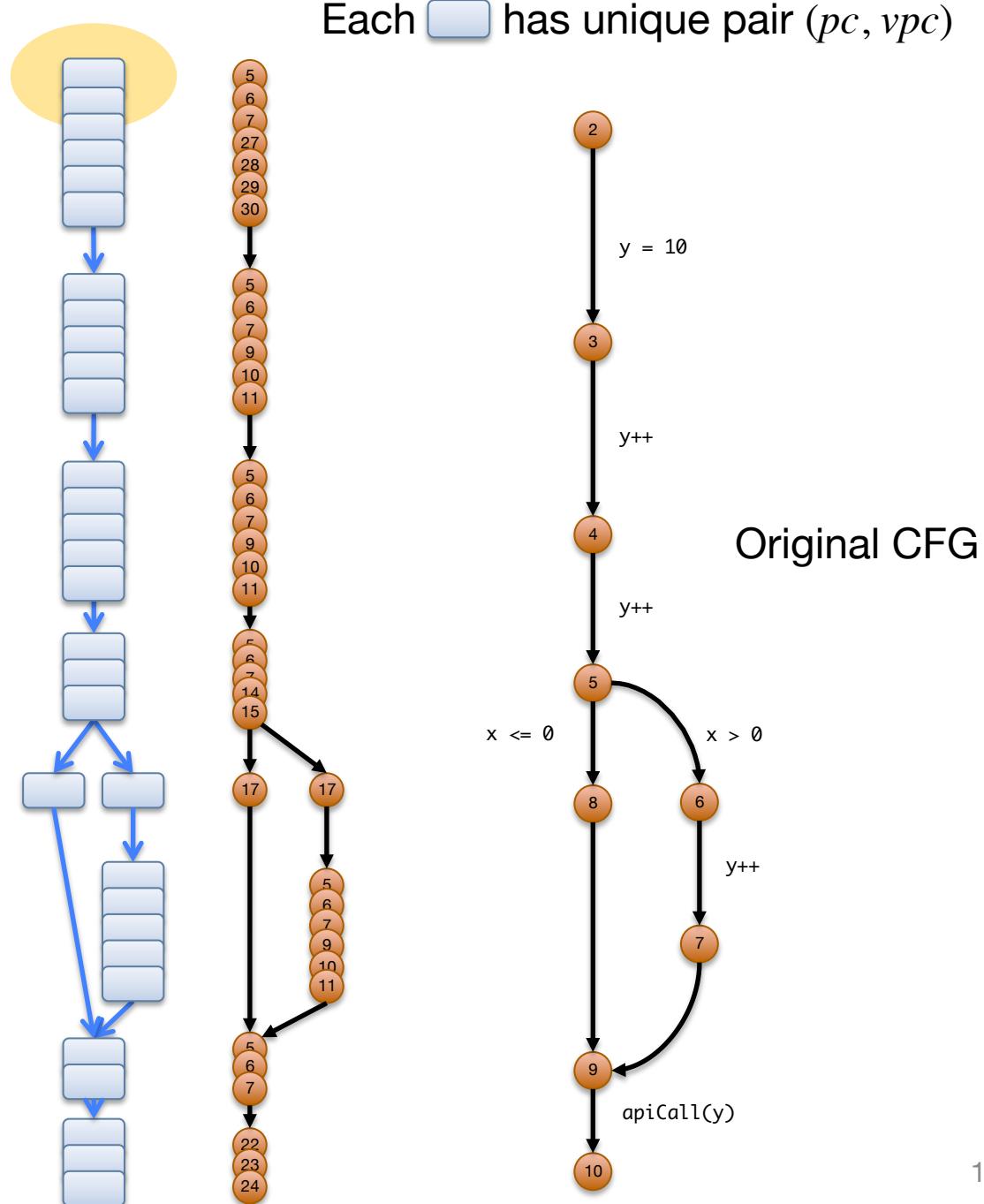
$$CFG = \{((\ell, v), (\ell', v')) \mid (\ell, v, a) \mapsto (\ell', v', a'), \\ \ell, \ell' \in L, a, a' \in A, v, v' \in V\}$$

Each  has unique pair (pc, vpc)



```

switch(code[vpc]) {
    case 03:
        op1 = code[vpc + 1];
        data[op1]++;
        vpc += 2;
        break;
    
```



- Constant propagation
- Dead code elimination
- Jump threading

Implementation

- Implemented in Jakstab [CAV'08]
 - Processes *obfuscated binaries*
 - Reconstructs CFGs in presence of indirect jumps
- Analysis
 - VPC-lifted Bounded Address Tracking [FMCAD'10]
 - Sets of memory states up to per-variable bound

$$\left\{ \left(\begin{array}{c} x = 5 \\ y = 20 \end{array} \right), \left(\begin{array}{c} x = 6 \\ y = 20 \end{array} \right), \left(\begin{array}{c} x = 7 \\ y = 24 \end{array} \right) \right\} \cup \left\{ \left(\begin{array}{c} x = 8 \\ y = 24 \end{array} \right), \left(\begin{array}{c} x = 9 \\ y = 28 \end{array} \right) \right\}$$

$k = 3$

$$\left\{ \left(\begin{array}{c} x = 5 \\ y = 20 \end{array} \right), \left(\begin{array}{c} x = 6 \\ y = 20 \end{array} \right), \left(\begin{array}{c} x = \top \\ y = 24 \end{array} \right), \left(\begin{array}{c} x = \top \\ y = 28 \end{array} \right) \right\}$$

VPC Discovery

- Intuition: VPC changes frequently
 - *Each instruction will have a separate VPC value*
 - *VPC is the busiest variable of the interpreter*
- Detect VPC on the fly
 - *The variable that hits the value bound per variable first is promoted to VPC*
 - *Heuristic – not guaranteed to identify it correctly*

Preliminary Results

Benchmark	Baseline	Similarity	Time
tamperproof guard	0%	81%	13s
search tree	0%	100%	312s
matrix multiply	0%	100%	311s
stuxnet	0%	87%	319s

- Targets
 - *Created with research obfuscator (C. Collberg)*
- Similarity
 - *Graph edit distance using basic block markers*
 - *Baseline: similarity of obfuscated code to original*

Conclusion

- Virtualization causes domain flattening
 - *Strips one level of location sensitivity*
- VPC-lifting reintroduces location sensitivity
- CFG reconstruction by tracing VPC values
- Ongoing work
 - *Improve VPC discovery*
 - *Apply to real world obfuscators*

<http://www.jakstab.org>