

# All You Ever Wanted to Know About Dynamic Taint Analysis & Forward Symbolic Execution (but might have been afraid to ask)

(Yes, we were trying to overflow the title length field  
on the submission server)

Edward J. Schwartz, **Thanassis Avgerinos**, David Brumley

***A Few Things*** You Need to Know About  
Dynamic Taint Analysis  
&  
Forward Symbolic Execution  
(but might have been afraid to ask)

Edward J. Schwartz, **Thanassis Avgerinos**, David Brumley

# The Root of All Evil

Humans write programs



This Talk:  
Computers Analyzing Programs Dynamically at  
Runtime

# Two Essential Runtime Analyses

Detect Exploits

[Costa2005,Crandall2005,  
Newsome2005,Suh2004]

Detect  
packing in malware  
[Bayer2009,Yin2007]

Dynamic Taint Analysis:  
What values are derived from user input?

Automated Test Case  
Generation

[Cadar2008,Godefroid2005,Sen2005]

Input Filter Generation  
[Costa2007,Brumley2008]

Forward Symbolic Execution:  
What input will make execution reach *this* line of code?

# Our Contributions

Computers Analyzing Programs  
Dynamically at Runtime

1: Turn English descriptions into an ***algorithm***

– Operational Semantics

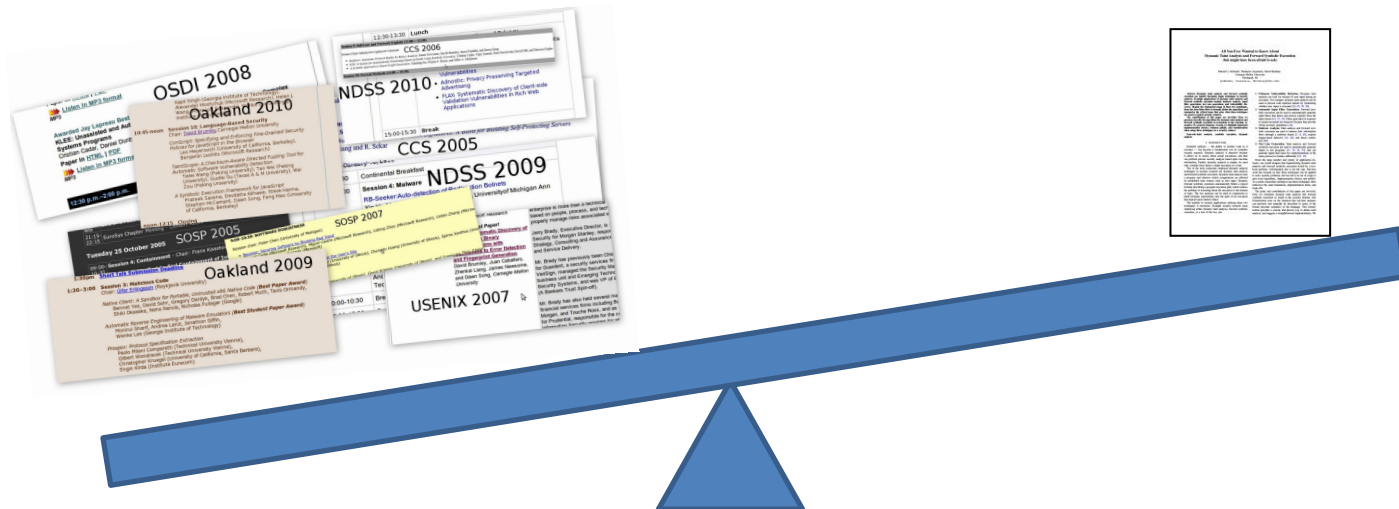
2: Algorithm highlights caveats, issues, and unsolved problems that are deceptively hard

Dynamic Taint Analysis:  
Is this value affected by user input?

Forward Symbolic Execution:  
What input will make execution reach ***this*** line of code?

# Our Contributions (cont'd)

3: Systematize recurring themes in a wealth of previous work



## Dynamic Taint Analysis: What values are derived from user input?

1. How it works – example
2. Desired properties
3. Example issue. Paper has many more.

● tainted      ● untainted

```

→ ● x = get_input(
    y = x + 42
    ...
    goto y
  
```



Input is tainted

Taint Introduction

$$\text{Input} \frac{t = \text{IsUntrusted}(src)}{\text{get\_input}(src) \downarrow t}$$

$\Delta$	
Var	Val
x	7

$\tau$	
Var	Tainted?
x	T



● tainted      ● untainted

●  $x = \text{get\_input}(\text{devil})$

→ ●  $y = \text{● } x + \text{● } 42$

...

goto y

Data derived from user input is tainted

$\Delta$	
Var	Val
x	7
y	49

## Taint Propagation

$$\text{BinOp} \frac{t_1 = \tau[x_1], t_2 = \tau[x_2]}{x_1 + x_2 \downarrow t_1 \vee t_2}$$

$\tau$	
Var	Tainted?
x	T
y	T

● tainted      ● untainted

●  $x = \text{get\_input}(\text{devil})$

●  $y = \text{● } x + \text{● } 42$

...



goto ●  $y$

Policy Violation Detected

Taint Checking

$P_{\text{goto}}(t_a) = \neg t_a$   
 (Must be true to execute)

$\Delta$	
Var	Val
$x$	7
$y$	49

$\tau$	
Var	Tainted?
$x$	T
$y$	T

# Difference Use: Exploit Detection

x = get\_input()

y = ...

...

goto y

Jumping to  
overwritten  
return address

```
...  
strcpy(buffer,argv[1]);  
...  
return ;
```

# Memory Load

Variables

$\Delta$

Var	Val
x	7

$\tau$

Var	Tainted?
x	T

Memory


$\mu$

Addr	Val
7	42

$\tau_\mu$

Addr	Tainted?
7	F

# Problem: Memory Addresses

→ **x** = get\_input()  
→ **y** = load(**x**)  
...  
goto y

All values derived  
from user input  
are tainted??

 $\Delta$ 

Var	Val
x	7

 $\mu$ 

Addr	Val
7	42

 $\tau_\mu$ 


Addr	Tainted?
7	F

# Policy 1: Taint depends only on the memory cell

```

x = get_in
y = load(
...
goto y

```



## Undertainting

Failing to identify tainted values  
- e.g., missing exploits

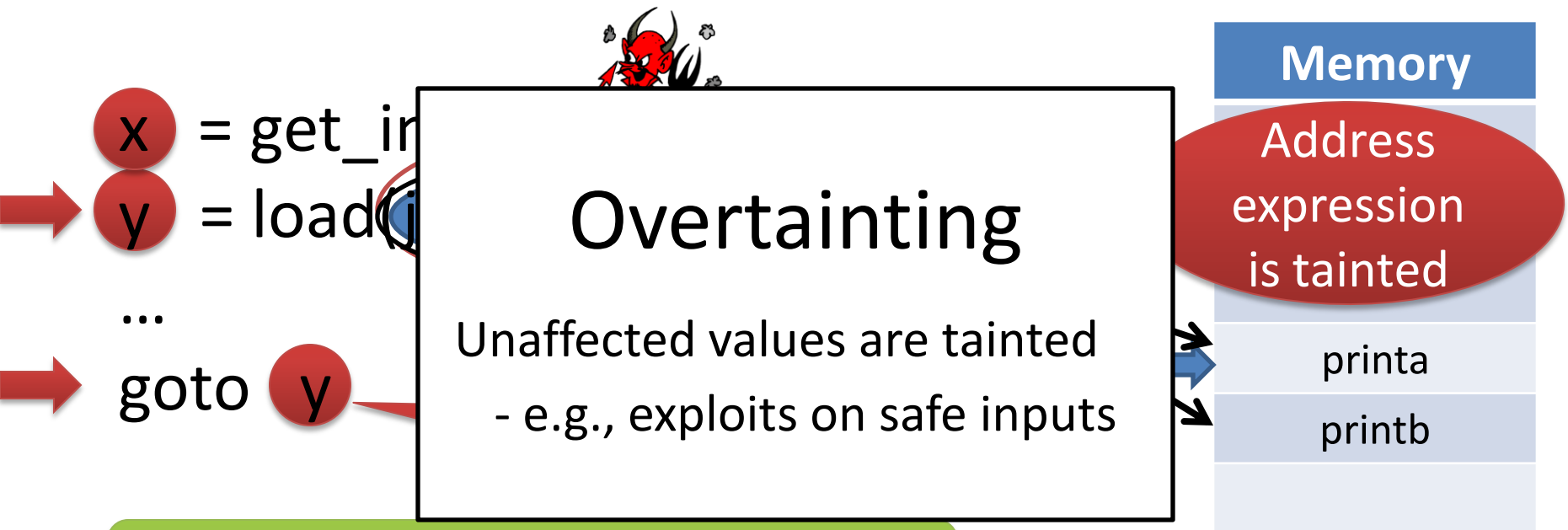
Var	Val
x	7
Addr	Val
7	42

## Taint Propagation

Load  $\frac{v = \Delta[x], t = \tau_\mu[v]}{\text{load}(x) \downarrow t}$

$\tau_\mu$ Addr	Tainted?
7	F

**Policy 2:** If either the address or the memory cell is tainted, then the value is tainted

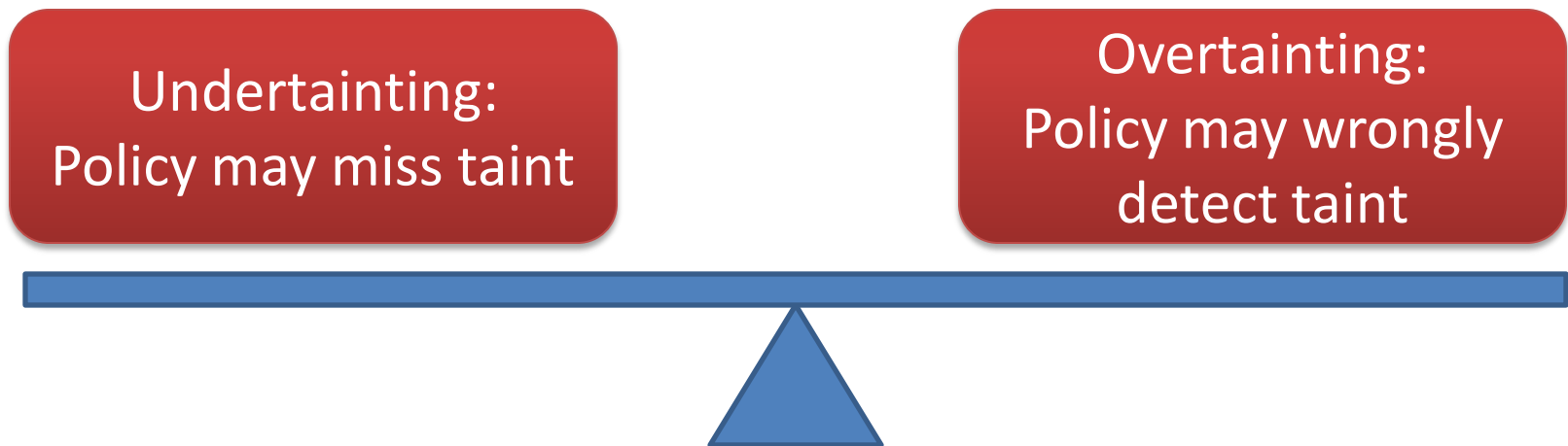


## Taint Propagation

Load  $\frac{v = \Delta[x], t = \tau_{\mu}[v], t_a = \tau[x]}{\text{load}(x) \downarrow t \ v \ t_a}$

# Research Challenge

State-of-the-Art is not perfect for all programs






Forward Symbolic Execution:  
What input will make execution reach *this* line of code?

- How it works – example
- Inherent problems of symbolic execution
- Proposed solutions

# The Challenge



$2^{32}$  possible  
inputs

0x12345678

```
packet_len(int header, char *packet)
char buf[2048] = "...";
if (header < 0)
    return 0;
if (header == 0x12345678)
    strcpy(buf, packet);
return strlen(buf);
```

Forward Symbolic Execution:  
What input will make execution  
reach *this* line of code?

# A Simple Example

*header symbolic can have any*

What input will make execution reach *this* line of code?

*header ≥ 0*

```
packetInterpreter(header, ...)
```

```
if (Interpreter())
```

```
return 0;
```

```
if (header < 0) Interpreter(0x12345678)
```

*header < 0*

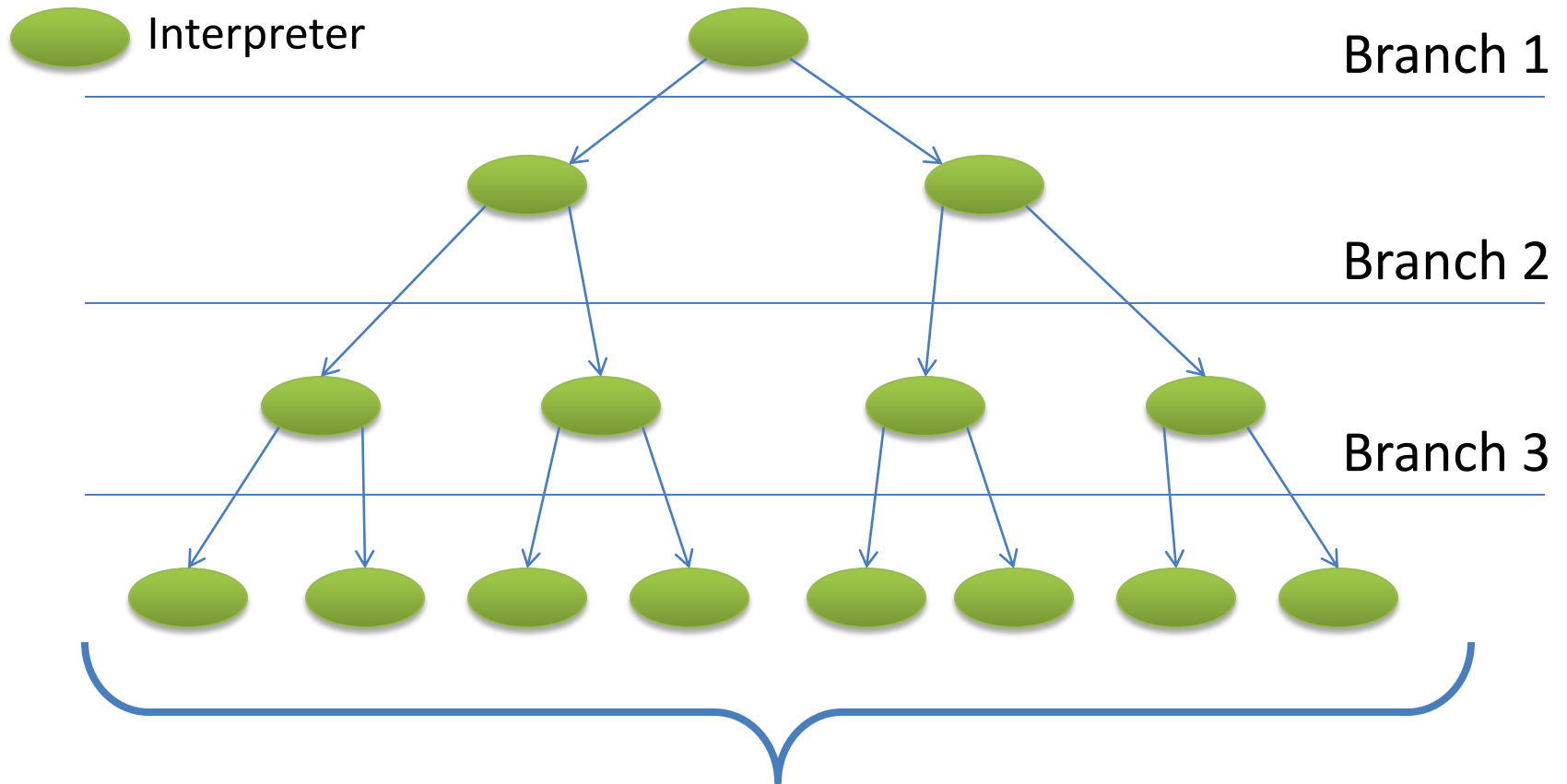
```
return strlen(buf);
```

```
strcpy(buf, packet);
```

*header ≥ 0 ∧ header != 0x12345678*

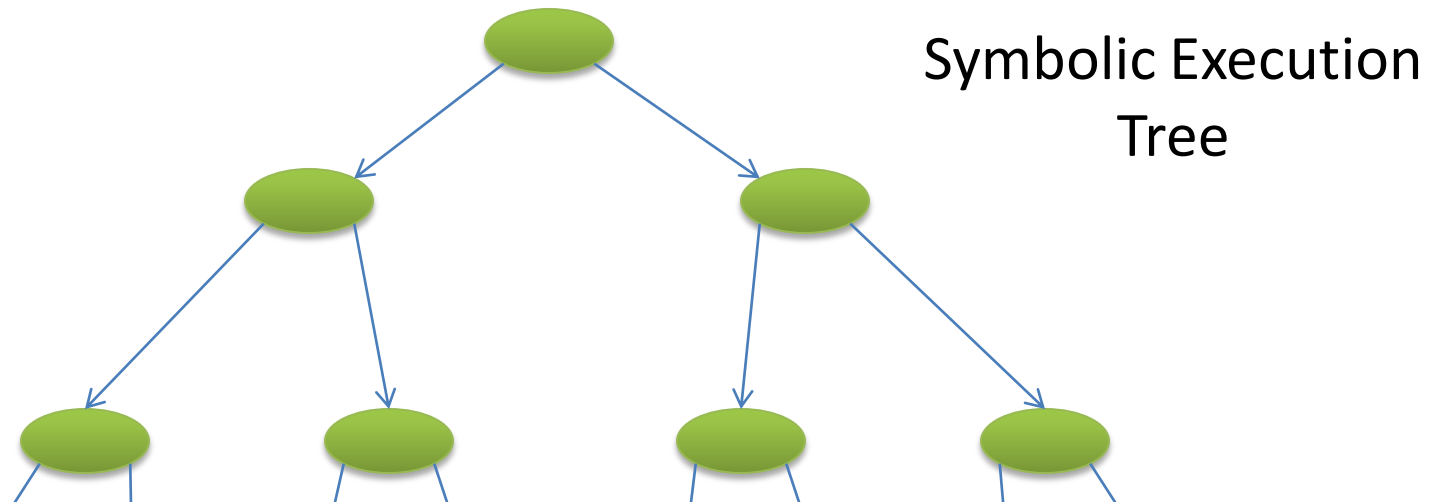
*header ≥ 0 ∧ header == 0x12345678*

# One Problem: Exponential Blowup Due to Branches



**Exponential Number of Interpreters/formulas in # of branches**

# Path Selection Heuristics



However, these are heuristics. In the worst case all create an exponential number of formulas in the tree height.

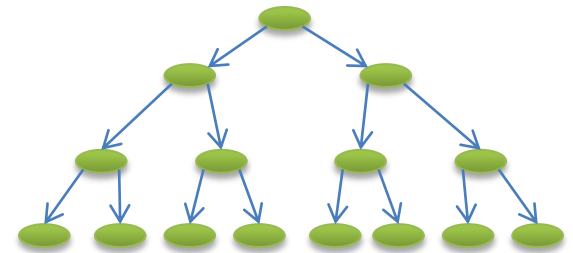
...

- Depth-First Search (bounded) , Random Search [Cadar2008]
- Concolic Testing [Sen2005,Godefroid2008]

# Symbolic Execution is *not* Easy

- Exponential number of interpreters/formulas

branching



- Exponentially-sized formulas

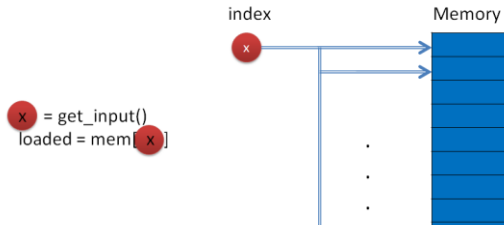
substitution



- Solving a formula is NP-Complete!

# Other *Important* Issues

## Symbolic Memory



## FORWARD SYMBOLIC EXECUTION

"Like normal execution, where inputs are substituted by symbolic variables"

King et al., 1976

## Formalization

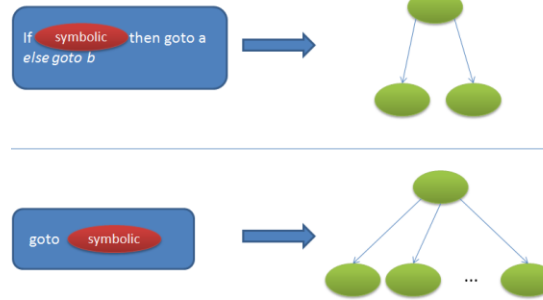
Other reasonable assumption?

```

e is input from src
T_INPUT  T_CONST
T_LOAD  T_STORE
T_ASSIGN T_STORE
T_COND  T_COND
T_GOTO  T_GOTO

```

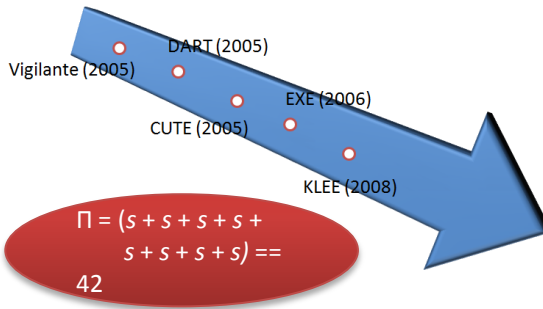
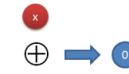
## Symbolic Jumps



## Sanitization



## More complex policies



# Conclusion

- Dynamic taint analysis and forward symbolic execution used extensively in literature
  - Formal algorithm and what is done for each possible step of execution often not emphasized
- We provided a formal definition and summarized
  - Critical issues
  - State-of-the-art solutions
  - Common tradeoffs



# Thank You!

thanassis@cmu.edu

# Questions?