

# Taint Analysis

# Paper

- Edward J. Schwartz, Thanassis Avgerinos, and David Brumley
  - All You Ever Wanted to Know about Dynamic Taint Analysis and Forward Symbolic Execution (but Might Have Been Afraid to Ask).
  - IEEE Symposium on Security and Privacy 2010

# 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?

 tainted
  untainted


  $x = \text{get\_input}(\text{devil})$   
 $y = x + 42$   
 ...  
 goto  $y$



Input is tainted

## Taint Introduction

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



$$\Delta$$





Var	Val
$x$	7

$$\tau$$

Var	Tainted?
$x$	T

 tainted
  untainted

  $x$  = get\_input(
 


  $y$  =   $x$  +  42

...

goto  $y$

Data derived from user input is tainted



## Taint Propagation

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

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

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

 tainted
  untainted

  $x$  = get\_input(
 

  $y$  =   $x$  +  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

## Real 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

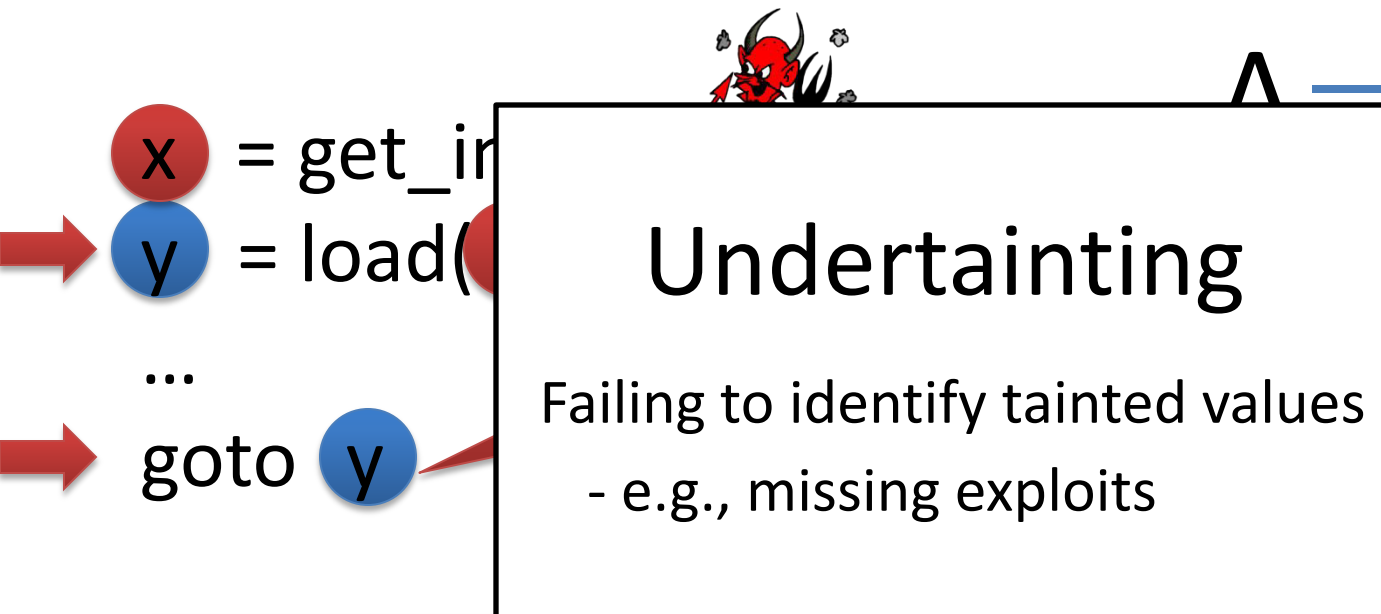
All values derived  
from user input  
are tainted??

$\Delta$	Var	Val
	x	7

$\mu$	Addr	Val
	7	42

$\tau_\mu$	Addr	Tainted?
	7	F

# Policy 1: Taint depends only on the memory cell



Var	Val
x	7
Addr	Val
7	42

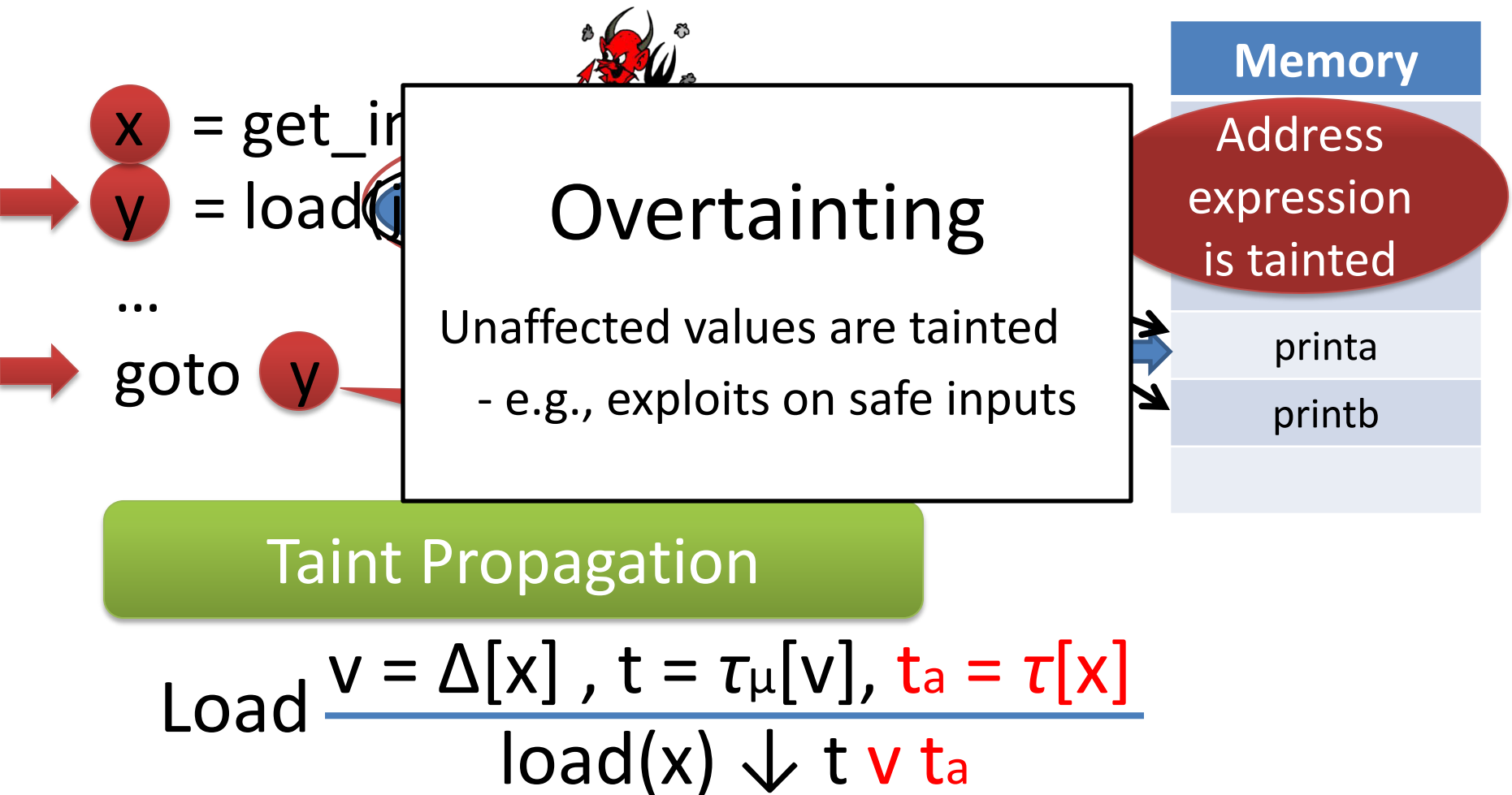
## Taint Propagation

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

$\tau_\mu$

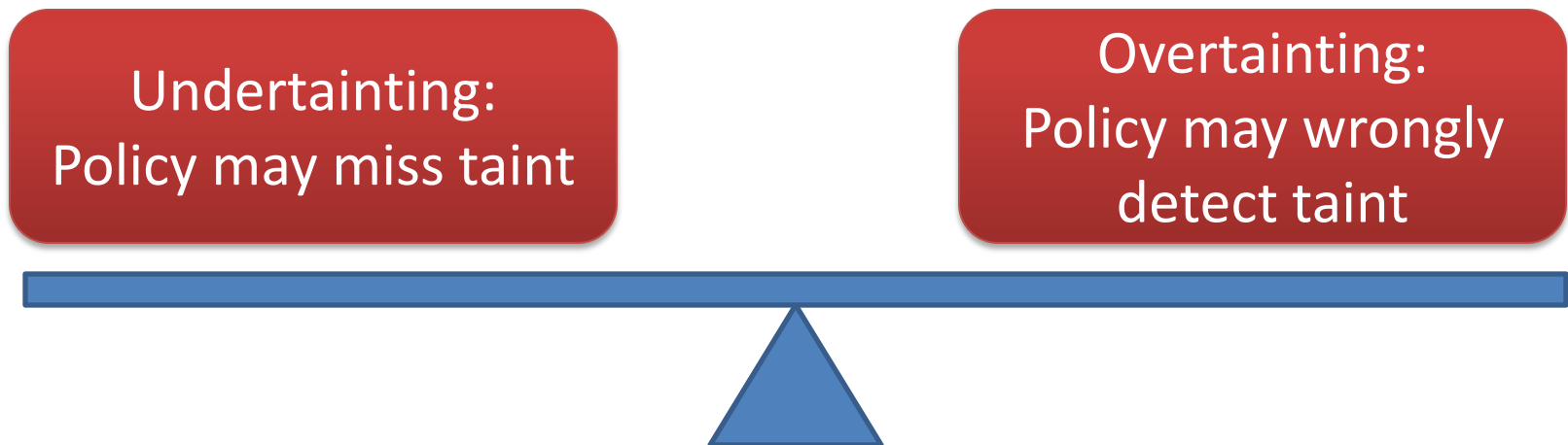
Addr	Tainted?
7	F

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




# Research Challenge

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



# The Challenge



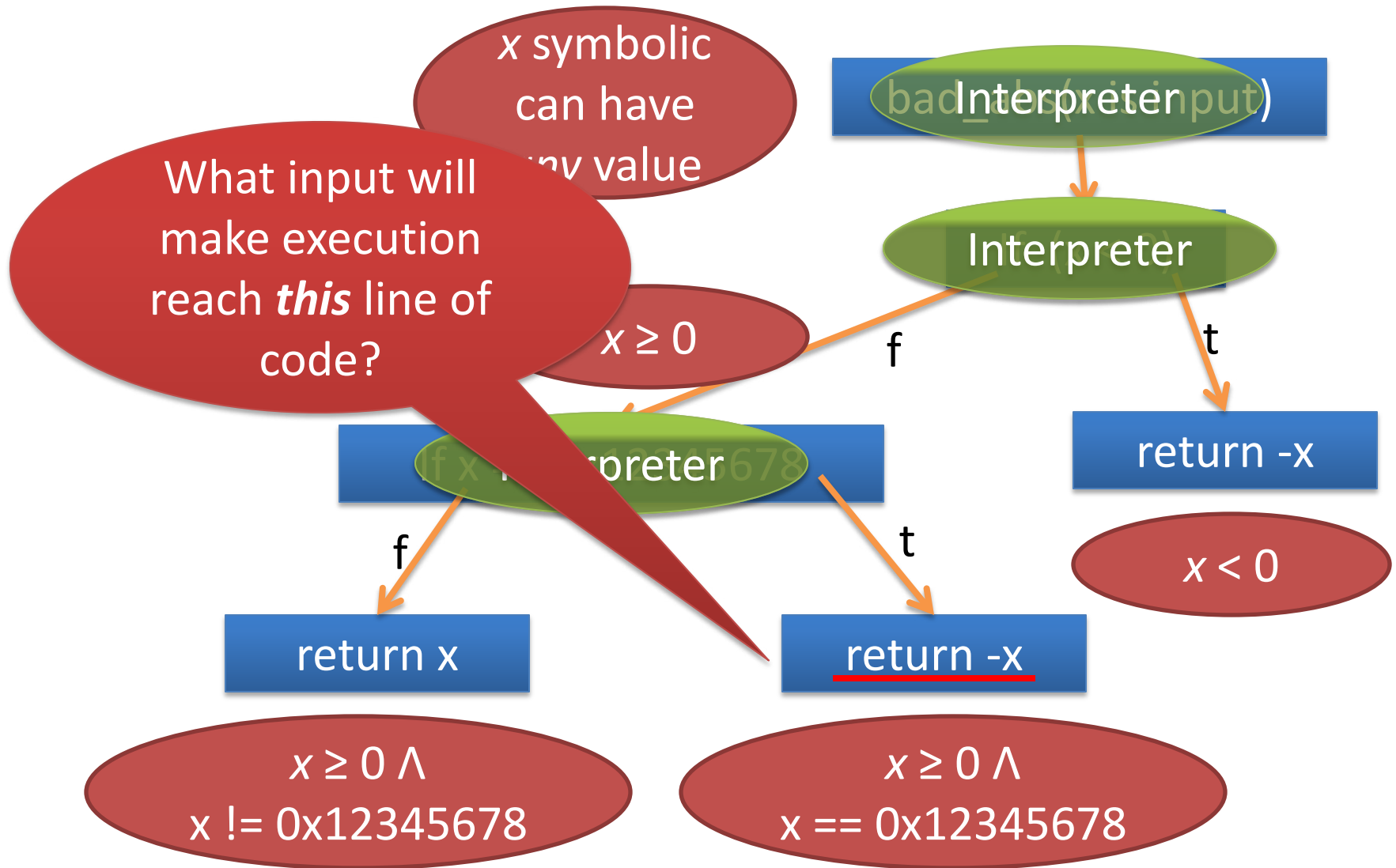
$2^{32}$  possible  
inputs

0x12345678

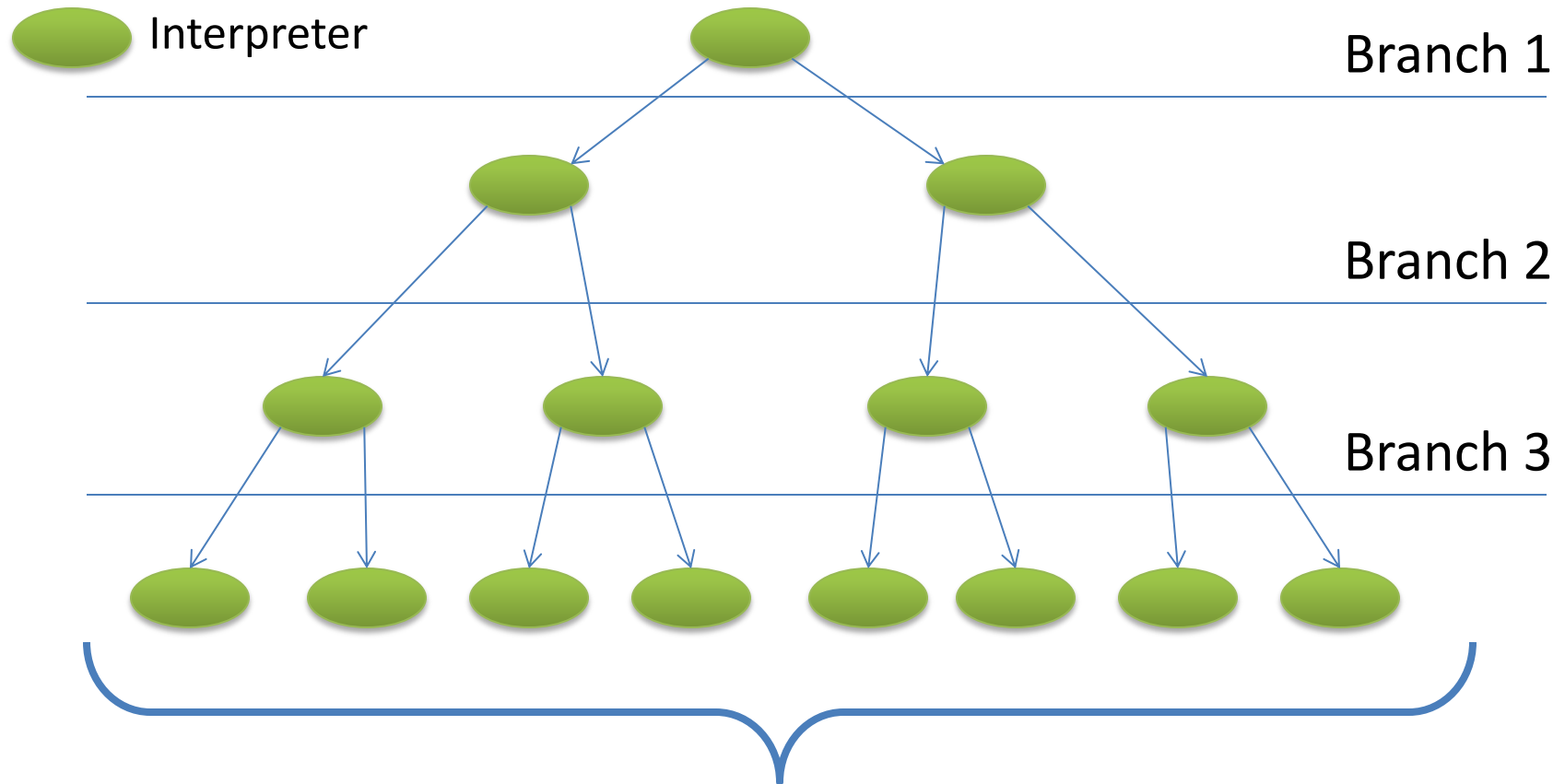
```
bad_abs(x is input)
  if (x < 0) then
    return -x
  if (x = 0x12345678) then
    return -x
  return x
```

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

# A Simple Example

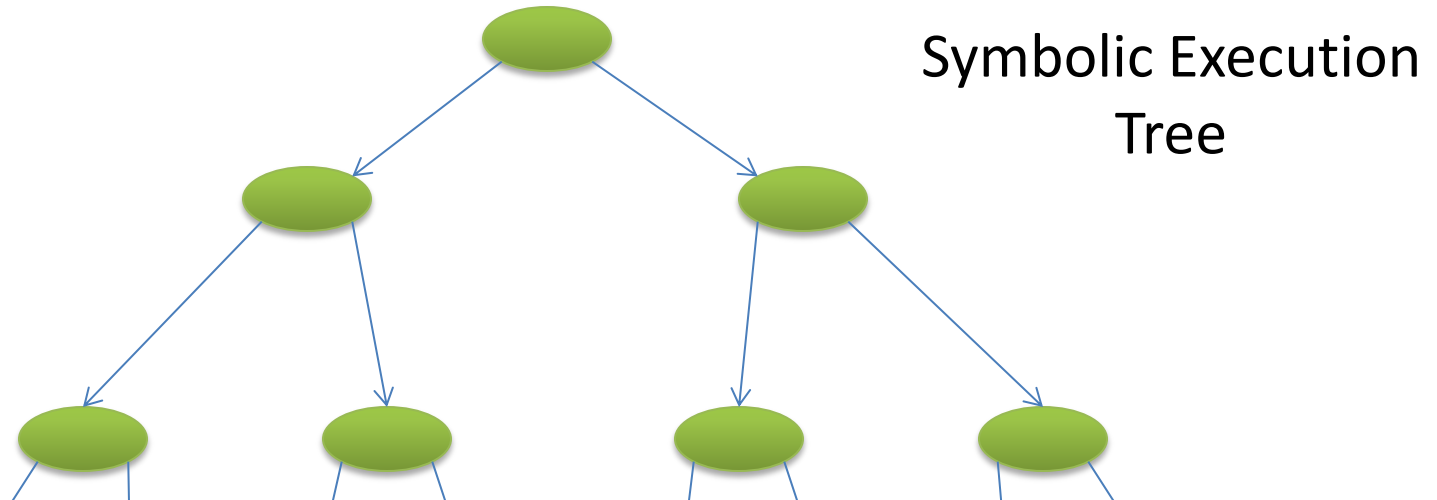


# One Problem: Exponential Blowup Due to Branches



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.

...

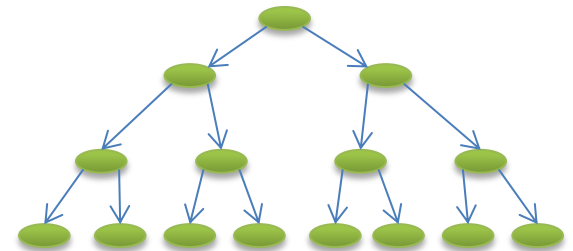
- Depth-First Search (bounded) ,Random Search [Cadar2008]
- Concolic Testing [Sen2005,Godefroid2008]



# Symbolic Execution is *not* Easy

- Exponential number of interpreters/formulas

branching



- Exponentially-sized formulas

substitution

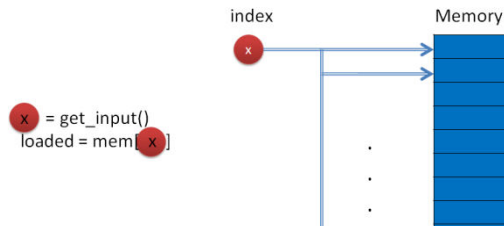


$s + s + s + s +$   
 $s + s + s + s == 42$

- Solving a formula is NP-Complete!

# Other *Important* Issues

## Symbolic Memory



## FORWARD SYMBOLIC EXECUTION

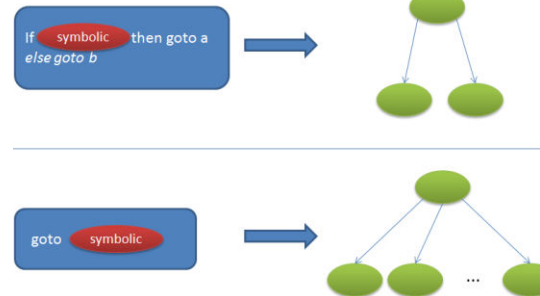
"Like normal execution, where inputs are substituted by symbolic variables"

King et al., 1976

## Formalization

$e$ is input from $src$	$T\_INPUT$	$\tau_p, \tau_a, \mu, \Delta \vdash \text{get\_input}() \rightarrow v \in \{v, \text{Pname}\}$	$T\_CONST$	$\tau_p, \tau_a, \mu, \Delta \vdash v \in \{v, \text{Pname}\}$
$\tau_p, \tau_a, \mu, \Delta \vdash \text{var} \in \{\Delta(\text{var}), \text{Pname}\}$	$T\_VAR$	$\tau_p, \tau_a, \mu, \Delta \vdash \text{var} \in \{v, \text{Pname}\}$	$T\_LOAD$	$\tau_p, \tau_a, \mu, \Delta \vdash \text{load } v \in \{v, \text{Pname}\}$
$\tau_p, \tau_a, \mu, \Delta \vdash e \in \{v, \text{Pname}\}$	$T\_UNOP$	$\tau_p, \tau_a, \mu, \Delta \vdash e \in \{v, \text{Pname}\}$	$T\_BINOP$	$\tau_p, \tau_a, \mu, \Delta \vdash e \in \{v, \text{Pname}\}$
$\tau_p, \tau_a, \mu, \Delta \vdash e_1 \in \{v, \text{Pname}\}$	$T\_BINOP$	$\tau_p, \tau_a, \mu, \Delta \vdash e_2 \in \{v, \text{Pname}\}$	$T\_ASSIGN$	$\tau_p, \tau_a, \mu, \Delta \vdash e \in \{v, \text{Pname}\}$
$\tau_p, \tau_a, \mu, \Delta \vdash e \in \{v, \text{Pname}\}$	$T\_ASSIGN$	$\tau_p, \tau_a, \mu, \Delta \vdash e \in \{v, \text{Pname}\}$	$T\_STORE$	$\tau_p, \tau_a, \mu, \Delta \vdash e \in \{v, \text{Pname}\}$
$\tau_p, \tau_a, \mu, \Delta \vdash e \in \{v, \text{Pname}\}$	$T\_STORE$	$\tau_p, \tau_a, \mu, \Delta \vdash e \in \{v, \text{Pname}\}$	$T\_COND$	$\tau_p, \tau_a, \mu, \Delta \vdash e \in \{v, \text{Pname}\}$
$\tau_p, \tau_a, \mu, \Delta \vdash e \in \{v, \text{Pname}\}$	$T\_COND$	$\tau_p, \tau_a, \mu, \Delta \vdash e \in \{v, \text{Pname}\}$	$T\_FCOND$	$\tau_p, \tau_a, \mu, \Delta \vdash e \in \{v, \text{Pname}\}$
$\tau_p, \tau_a, \mu, \Delta \vdash e \in \{v, \text{Pname}\}$	$T\_FCOND$	$\tau_p, \tau_a, \mu, \Delta \vdash e \in \{v, \text{Pname}\}$	$T\_GOTO$	$\tau_p, \tau_a, \mu, \Delta \vdash e \in \{v, \text{Pname}\}$
$\tau_p, \tau_a, \mu, \Delta \vdash e \in \{v, \text{Pname}\}$	$T\_GOTO$	$\tau_p, \tau_a, \mu, \Delta \vdash e \in \{v, \text{Pname}\}$	$T\_GOOD$	$\tau_p, \tau_a, \mu, \Delta \vdash e \in \{v, \text{Pname}\}$

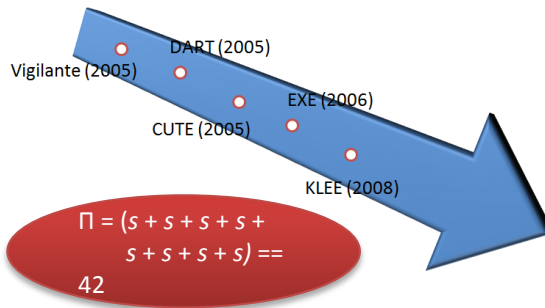
## Symbolic Jumps



## Sanitization



## More complex policies



# Summary

- 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