## Slide 1

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

5/19/2010     Carnegie Mellon University     1

## Slide 2

A *Few Things* You Need to Know About
Dynamic Taint Analysis
&
Forward Symbolic Execution
(but might have been afraid to ask)
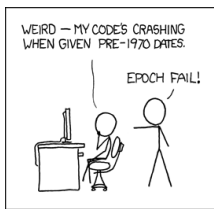
Edward J. Schwartz, Thanassis Avgerinos, David Brumley

5/19/2010     Carnegie Mellon University     2

## Slide 3
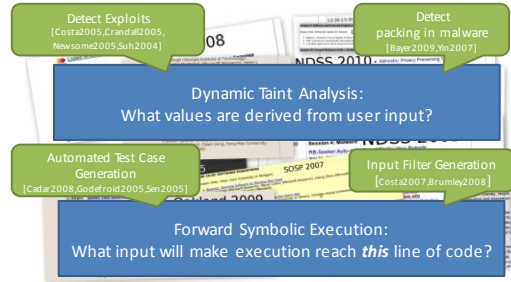
### The Root of All Evil

Humans write programs



WEIRD — MY CODE'S CRASHING
WHEN GIVEN PRE-1970 DATES.

EPOCH FAIL!

This Talk:
Computers Analyzing Programs Dynamically at Runtime

5/19/2010     Carnegie Mellon University     3

## Slide 4

### 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?

5/19/2010     Carnegie Mellon University     4

## Slide 5

### Our Contributions

Computers Analyzing Programs Dynamically at Runtime

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

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

1: Turn English descriptions into an *algorithm*
   – Operational Semantics
2: Algorithm highlights caveats, issues, and unsolved problems that are deceptively hard

5/19/2010     Carnegie Mellon University     5

## Slide 6

### Our Contributions (cont'd)

3: Systematize recurring themes in a wealth of previous work



5/19/2010     Carnegie Mellon University     6

**Slide 7**

| 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.

5/19/2010 — Carnegie Mellon University — 7

---

**Slide 8**

● tainted   ● untainted

$x$ = get_input( )
$y$ = $x$ + 42
...
goto $y$

Input is tainted

**Taint Introduction**

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

Δ

| Var | Val |
| --- | --- |
| x | 7 |

$\tau$

| Var | Tainted? |
| --- | --- |
| x | T |

5/19/2010 — Carnegie Mellon University — 8

---

**Slide 9**

● tainted   ● untainted

$x$ = get_input( )
$y$ = $x$ + 42
...
goto $y$

Data derived from user input is tainted

**Taint Propagation**

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

Δ

| Var | Val |
| --- | --- |
| x | 7 |
| y | 49 |

$\tau$

| Var | Tainted? |
| --- | --- |
| x | T |
| y | T |

5/19/2010 — Carnegie Mellon University — 9

---

**Slide 10**

● 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)

Δ

| Var | Val |
| --- | --- |
| x | 7 |
| y | 49 |

$\tau$

| Var | Tainted? |
| --- | --- |
| x | T |
| y | T |

5/19/2010 — Carnegie Mellon University — 10

---

**Slide 11**

$x$ = get_input( )
$y$ = ...
...
goto $y$

Jumping to overwritten return address

Real Use:
Exploit Detection

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

5/19/2010 — Carnegie Mellon University — 11

---

**Slide 12**

## Memory Load

Variables

Δ

| Var | Val |
| --- | --- |
| x | 7 |

$\tau$

| Var | Tainted? |
| --- | --- |
| x | T |

Memory

μ

| Addr | Val |
| --- | --- |
| 7 | 42 |

$\tau_\mu$

| Addr | Tainted? |
| --- | --- |
| 7 | F |

5/19/2010 — Carnegie Mellon University — 12

## Problem: Memory Addresses

```
x = get_input(        )
y = load( x )
...
goto  y
```

All values derived from user input are tainted??

Δ

| Var | Val |
|-----|-----|
| x | 7 |

μ

| Addr | Val |
|------|-----|
| 7 | 42 |

$\tau_\mu$

| Addr | Tainted? |
|------|----------|
| 7 | F |

5/19/2010     Carnegie Mellon University     13

## Policy 1: Taint depends only on the memory cell

```
x = get_i
y = load(
...
goto  y
```

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

Δ

| Var | Val |
|-----|-----|
| x | 7 |

| Addr | Val |
|------|-----|
| 7 | 42 |

**Taint Propagation**

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

$\tau_\mu$

| Addr | Tainted? |
|------|----------|
| 7 | F |

5/19/2010     Carnegie Mellon University     14

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

```
x = get_i
y = load(
...
goto  y
```

### Overtainting
Unaffected values are tainted
- e.g., exploits on safe inputs

**Memory**

Address expression is tainted
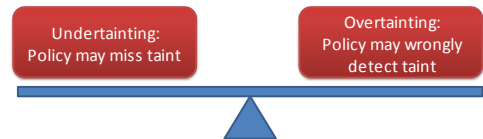
printa
printb

**Taint Propagation**

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

5/19/2010     Carnegie Mellon University     15

## Research Challenge
## State-of-the-Art is not perfect for all programs

Undertainting:
Policy may miss taint

Overtainting:
Policy may wrongly detect taint

5/19/2010     Carnegie Mellon University     16

---

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

- How it works – example

- Inherent problems of symbolic execution

- Proposed solutions

5/19/2010     Carnegie Mellon University     17

## The Challenge

2³² possible inputs

0x12345678

```
bad_abs(x is input)
    if (x < 0) then
        return -x
    if (x = 0x12345678) then
        return -x
    return x
```
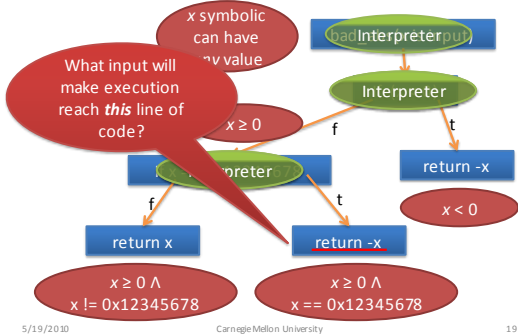
Forward Symbolic Execution:
What input will make execution reach **this** line of code?

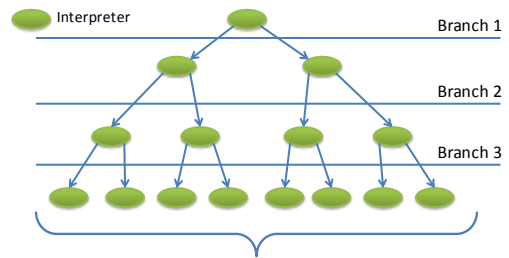5/19/2010     Carnegie Mellon University     18

## A Simple Example



x symbolic can have any value

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

bad Interpreter(put)

Interpreter

$x \geq 0$

f        t

Interpreter        return -x

f        t        $x < 0$

return x        return -x

$x \geq 0 \land$ x != 0x12345678

$x \geq 0 \land$ x == 0x12345678

## One Problem: Exponential Blowup Due to Branches



Interpreter        Branch 1

Branch 2

Branch 3

Exponential Number of Interpreters/formulas in # of branches

## Path Selection Heuristics



Symbolic Execution Tree

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]
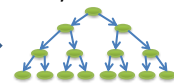
## Symbolic Execution is *not* Easy

- Exponential number of interpreters/formulas

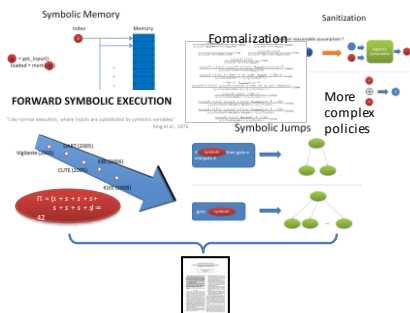  branching →



- Exponentially-sized formulas

  substitution →

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

- Solving a formula is NP-Complete!

## Other *Important* Issues

## 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?