FINDING INFORMATION FLOW BUGS WITH SYMBOLIC EXECUTION

JOSÉ FRAGOSO SANTOS ASSISTANT PROFESSOR, DEI

NOVEMBER 2019 @ SOFTWARE SECURITY (MEIC)

CONTENTS

1. Motivation: Information flow bugs on the Web

GOAL: UNDERSTAND HOW TO USE SELF-COMPOSITION TO FIND INFORMATION FLOW BUGS

4. Concolic Symbolic Execution

RESEARCH PROGRAM – VERIFICATION AND TESTING

Verification and Testing of JavaScript (Web) Programs



RESEARCH PROGRAM – VERIFICATION AND TESTING

Verification and Testing of JavaScript (Web) Programs



RESEARCH PROGRAM – VERIFICATION AND TESTING

Verification and Testing of JavaScript (Web) Programs



RESEARCH PROGRAM – INFORMATION FLOW CONTROL

Information Flow Control for Web Programs



I. INFORMATION FLOW CONTROL ON THE WEB

JS Security Vulnerabilities – Browser Extensions

Browser extensions are often implemented directly in JavaScript

- Browser extensions execute with elevated privileges
- Web apps can communicate with the extensions executing in the browser via the Browser API
- Malicious Web Apps can exploit browser extension leaks to obtain user-sensitive data

EmPoWeb* found often than **100** leaks in JS browser extensions

*EmPoWeb: Empowering Web Applications with Browser Extensions, Dolière Francis Somé, 2019

JS SECURITY VULNERABILITIES – BROWSER EXTENSIONS

EmPoWeb* found often than **100** leaks in JS browser extensions

SIMPLE SYNTACTIC ANALYSIS ARE NOT ENOUGH!

chrome.cookies.remove

But there are multiple other ways to access these values that would not be caught by this analysis

Leaking 1 bit

```
function getBit(x) {
    var a = [ "a", "b", "c" ];
    if ((x % 2) === 0) {
        Object.defineProperty(a, "1", {value: "d", configurable: false})
    }
    a.length = 1;
    return (a.length === 1) ? 1 : 0;
}
```

Implicit Information flow via a property descriptor

Leaking 1 bit

```
function getBit(x) {
    var op = {};
    var f = function () {};
    f.prototype = op;
    var o = new f ();
    if ((x % 2) === 0) {
        Object.defineProperty(op, "foo", {value: 0, configurable: false})
    }
    o.foo = 1;
    return o.foo;
}
```

Implicit Information flow via prototype inheritance

var x = chrome;

Bit by bit, we can leak all bits

```
function getBits(x) {
    var bits = [];
    while (x > 0) {
        bits.push(getBit(x));
        x = x >> 1;
    }
    return (bits.length === 0) ? 0 : bits.reverse();
}
```

And now we can learn:

chrome.cookies.remove

```
var x = chrome;
var y1 = "cook";
var y2 = "ies";
var z1 = "rem";
var z2 = "ove";
var x = getBits(x[y1+y2][z1+z2])
```



Take-home message

Malicious code can exploit corner case behaviors of the JavaScript semantics to encode sophisticated information flow

We need to do better than a simple syntactic analysis!

II. PROGRAM PROPERTIES

VERIFICATION VS BUG-FINDING

Program property = set of "behaviors"

Verification: verifying a given program *P* with respect to a given property *S* means proving that all the behaviors of *P* are contained in *S*

$\llbracket P \rrbracket \subseteq \mathcal{S}$

VERIFICATION VS BUG-FINDING

Program property = set of "behaviors"

Bug Finding: debugging a program *P* with respect to a given property *S* means finding a behavior of *P* that is not in *S*

$$\llbracket P \rrbracket \cap \overline{\mathcal{S}} \neq \emptyset$$

VERIFICATION VS BUG-FINDING

Verification Bug-finding

 $\llbracket P \rrbracket \subseteq \mathcal{S} \qquad \llbracket P \rrbracket \cap \overline{\mathcal{S}} \neq \emptyset$

Hard: for/all

Easy: exists

PROGRAM PROPERTIES

Question: How do we define the set of behaviors?

$\llbracket P \rrbracket$?

Depending on how we define the set of allowed behaviors, we get different classes of properties.

TRACE PROPERTIES



Leslie Lamport

Safety Properties: Nothing bad ever happens

Liveness Properties: Something good eventually happens

TRACE PROPERTIES

Safety Properties: Nothing bad ever happens

- Type Safety
- Memory Safety (no null pointer exceptions)

Liveness Properties: Something good eventually happens

TRACE PROPERTIES

Safety Properties: Nothing bad ever happens

- Type Safety
- Memory Safety (no null pointer exceptions)

Liveness Properties: Something good eventually happens

- Termination
- Absence of memory leaks

Program property = set of "behaviors"

How do we formally define program traces?

Trace property = set of program traces

OPERATIONAL SEMANTICS - WHILE LANGUAGE

Syntax of While

$$e_1, e_2 \in \mathcal{E} \triangleq n \mid x \mid \ominus e_1 \mid e_1 \oplus e_2$$

$$s_1, s_2 \in \mathcal{S} \triangleq \text{skip} \mid x := e \mid s_1; s_2$$

$$\mid \text{if}(e) \{ s_1 \} \text{else} \{ s_2 \}$$

$$\mid \text{while}(e) \{ s_1 \}$$

Small-Step Operational Semantics



Gordon Plotkin

OPERATIONAL SEMANTICS - WHILE LANGUAGE

Syntax of While

$$e_1, e_2 \in \mathcal{E} \triangleq n \mid x \mid \ominus e_1 \mid e_1 \oplus e_2$$
$$s_1, s_2 \in \mathcal{S} \triangleq \text{skip} \mid x := e \mid s_1; s_2$$
$$\mid \text{if}(e) \{ s_1 \} \text{else} \{ s_2 \}$$
$$\mid \text{while}(e) \{ s_1 \}$$

Small-step Transition

 $\langle \rho, s \rangle \rightarrow \langle \rho', s' \rangle$

State = Variable Store

 $\rho \in \text{Store} : \text{PVar} \rightarrow \mathbb{N}$

A SIMPLE WHILE LANGUAGE - SEMANTICS

Assignment	IF - TRUE
$\rho' = \rho \left[x \mapsto \llbracket e \rrbracket_{\rho} \right]$	$[\![e]\!]_{\rho} \neq 0$
$\overline{\langle \rho, x := e \rangle} \to \langle \rho', \operatorname{skip} \rangle$	$\overline{\langle \rho, if(e) \{ s_1 \}}$ else $\{ s_2 \} \rangle \rightarrow \langle \rho, s_1 \rangle$
If - False	Seq - 1
$[\![e]\!]_{\rho} = 0$	$\langle \rho, s_1 \rangle \to \langle \rho, s_1' \rangle$
$\langle \rho, \text{if}(e) \{ s_1 \} \text{else} \{ s_2 \} \rangle$	$\rightarrow \langle \rho, s_2 \rangle \qquad \overline{\langle \rho, s_1; s_2 \rangle} \rightarrow \langle \rho, s_1'; s_2 \rangle$
SEQ - 2 $(a, a) \rightarrow (a, a)$	WHILE $s' = if(e)\{s; while(e)\{s\}\}else\{skip\}$
$\langle p, s_1 p; s_2 \rangle \rightarrow \langle p, s_2 \rangle$	$\langle \rho, \text{while}(e)\{s\} \rangle \rightarrow \langle \rho, s' \rangle$

A SIMPLE WHILE LANGUAGE - SEMANTICS

Division by 0 generates $\frac{1}{2}$

Assignment - Error	If - Error
$\llbracket e \rrbracket_{\rho} = \cancel{2}$	$\llbracket e \rrbracket_{\rho} = \measuredangle$
$\langle \rho, x := e \rangle \rightarrow \langle \sharp, \operatorname{skip} \rangle$	$\langle \rho, if(e) \{ s_1 \} else \{ s_2 \} \rangle \rightarrow \langle 2, skip \rangle$

Trace property = set of program traces

Program Trace =?

 $\llbracket s \rrbracket \triangleq \{ [\langle \rho_0, s_0 \rangle, ..., \langle \rho_n, s_n \rangle] \mid s_0 = s \land s_n = \text{skip} \land \forall_{0 \le i < n} \langle \rho_i, s_i \rangle \to \langle \rho_{i+1}, s_{i+1} \rangle \}$

Trace property = No division by 0

$$\begin{aligned} \text{NoDivZero} &\triangleq \{ [\langle \rho_0, s_0 \rangle, ..., \langle \rho_n, s_n \rangle] \mid s_n = \text{skip} \land \rho_n \neq \notin \\ &\land \forall_{0 \leq i < n} \langle \rho_i, s_i \rangle \to \langle \rho_{i+1}, s_{i+1} \rangle \end{aligned}$$

Trace property = Termination

$$\begin{array}{l} \texttt{Termination} \triangleq \{ [\langle \rho_0, s_0 \rangle, ..., \langle \rho_n, s_n \rangle] \mid s_n = \texttt{skip} \\ & \wedge \forall_{0 \leq i < n} \langle \rho_i, s_i \rangle \rightarrow \langle \rho_{i+1}, s_{i+1} \rangle \} \end{array}$$

Tx0 = Programs that terminate with **x** set to **0**

T0 = Programs that terminate with a **all variables** set to **0**

TRACE PROPERTIES - SUMMARY



WHAT ABOUT NON-INTERFERENCE?



J. Meseguer

J. Goguen

$\mathcal{NI}(\Gamma) \triangleq \{s \mid \forall \rho_1, \rho_2 \, . \, \rho_1 =_L^{\Gamma} \rho_2 \, \land \, \langle \rho_1, s \rangle \to^* \langle \rho'_1, \mathsf{skip} \rangle \\ \wedge \langle \rho_2, s \rangle \to^* \langle \rho'_2, \mathsf{skip} \rangle \implies \rho'_1 =_L^{\Gamma} \rho'_2 \}$

Security Labeling

 $\Gamma: \mathsf{PVar} \to \mathcal{L} \qquad \mathcal{L} = \langle \{L, H\}, \sqsubseteq \rangle$

Security Labeling

NON-INTERFERENCE: 2-TRACE PROPERTY

Non-Interference is a **2-trace** property

$$\mathcal{NI}(\Gamma) \triangleq \{s \mid \forall \rho_1, \rho_2 . \rho_1 =_L^{\Gamma} \rho_2 \land \langle \rho_1, s \rangle \to^* \langle \rho'_1, \text{skip} \rangle \\ \land \langle \rho_2, s \rangle \to^* \langle \rho'_2, \text{skip} \rangle \implies \rho'_1 =_L^{\Gamma} \rho'_2 \}$$

$$\begin{split} \mathcal{NI}(\Gamma) &\triangleq \{ ([\langle \rho_0, s_0 \rangle, ..., \langle \rho_n, s_n \rangle], [\langle \rho'_0, s'_0 \rangle, ..., \langle \rho'_m, s'_m \rangle]) \\ &\quad | \forall_{0 \leq i < n} \langle \rho_i, s_i \rangle \rightarrow \langle \rho_{i+1}, s_{i+1} \rangle \\ &\quad \wedge \forall_{0 \leq i < m} \langle \rho'_i, s'_i \rangle \rightarrow \langle \rho'_{i+1}, s'_{i+1} \rangle \\ &\quad \wedge s_0 = s'_0 \wedge (\rho_0 = {}^{\Gamma}_L \rho'_0 \Longrightarrow \rho_n = {}^{\Gamma}_L \rho'_m) \} \end{split}$$

INFORMATION FLOW BUG

$$\mathcal{NI}(\Gamma) \triangleq \{s \mid \forall \rho_1, \rho_2 \, . \, \rho_1 =_L^{\Gamma} \rho_2 \, \land \, \langle \rho_1, s \rangle \to^* \langle \rho'_1, \mathsf{skip} \rangle \\ \wedge \langle \rho_2, s \rangle \to^* \langle \rho'_2, \mathsf{skip} \rangle \implies \rho'_1 =_L^{\Gamma} \rho'_2 \}$$

A pair of stores (ρ_1, ρ_2) that prove that:

$$s \notin \mathcal{NI}(\Gamma)$$

$$\langle \rho_1, s \rangle \to^* \langle \rho'_1, \text{skip} \rangle \land \langle \rho_2, s \rangle \to^* \langle \rho'_2, \text{skip} \rangle$$

$$\wedge \rho_1 = {}_L^{\Gamma} \rho_2 \land \rho'_1 \neq {}_L^{\Gamma} \rho'_2$$

Hyper-Properties



F. Schneider

2-Trace Properties: Properties of 2 traces

- Non-Interference
- Dependency

N-Trace Properties: Meta-dependencies
PROGRAM PROPERTIES - SUMMARY



II. SELF-COMPOSITION + SYMBOLIC EXECUTION

Idea: Reduce non-interference to a **safety property** by transpiling the given program

$s \in \mathcal{NI}(\Gamma) \iff \llbracket \mathcal{C}(s) \rrbracket \subseteq \mathcal{T}(\Gamma)$

 $\mathcal{T}(\Gamma)$ - a safety property that only depends on Γ

C - a transpiler that computes the self-composition of s $T_{. Rezk}$

SELF-COMPOSITION - THE MAIN IDEA







P. D'Argenio



SELF-COMPOSITION – WHY?

Idea trar DESIGN AND IMPLEMENT, ESPECIALLY WHEN TARGETING REAL-WORLD LANGUAGES

Wh safe WE ARE GOING TO USE SYMBOLIC EXECUTION scra



Does the **assertion** hold?



We are going to execute the generated program symbolically

Instead of using concrete values, we use **symbolic variables**

assume(l1 = l2); l1 := h1; l2 := h2; assert(l1 = l2)

1 := h

True, [l1→#l1, h1→#h1, l2→#l2, h2→#h2]
Symbolic Store

Path Condition: conjunction of all the expressions on which the execution has branched before reaching the current execution point

1 := h True, [$l1 \rightarrow \#l1$, $h1 \rightarrow \#h1$, $l2 \rightarrow \#l2$, $h2 \rightarrow \#h2$] assume(11 = 12) #11=#12, [$11\rightarrow\#11$, $h1\rightarrow\#h1$, $12\rightarrow\#12$, $h2\rightarrow\#h2$] l1 := h1 assume(11 = 12);11 := h1;#11=#12, [$11\rightarrow\#h1$, $h1\rightarrow\#h1$, $12\rightarrow\#12$, $h2\rightarrow\#h2$] 12 := h2;12 := h2 assert(l1 = l2)#11=#12, [$11\rightarrow\#h1$, $h1\rightarrow\#h1$, $12\rightarrow\#h2$, $h2\rightarrow\#h2$]

1 := h 12 := h2 #11=#12, [$11\rightarrow\#h1$, $h1\rightarrow\#h1$, $12\rightarrow\#h2$, $h2\rightarrow\#h2$] assert(11 = 12)assume(11 = 12);11 := h1; $(\#11 = \#12) \Rightarrow (\#h1 = \#h2)$ Valid? 12 := h2;assert(l1 = l2) $(\#11 = \#12) \land (\#h1 \neq \#h2)$ **SAT?**

1 := h 12 := h2 #11=#12, [$11\rightarrow\#h1$, $h1\rightarrow\#h1$, $12\rightarrow\#h2$, $h2\rightarrow\#h2$] assert(11 = 12)assume(11 = 12);11 := h1; $(\#11 = \#12) \land (\#h1 \neq \#h2)$ SAT? 12 := h2;assert(l1 = l2)**YES!** [#l1→0, #h1→0, #l2→0, #h2→1]



True, [
$$l1 \rightarrow \#l1$$
, $h1 \rightarrow \#h1$, $l2 \rightarrow \#l2$, $h2 \rightarrow \#h2$]
assume($l1 = l2$)
$l1=\#l2$, [$l1 \rightarrow \#l1$, $h1 \rightarrow \#h1$, $l2 \rightarrow \#l2$, $h2 \rightarrow \#h2$]
if ($h1$)
$h1 = 0$
$h1 \neq 0$
Next Slide

assume(l1 = l2);
if (h1) {
 l1 := 1
} else { skip };
if (h2) {
 l2 := 1
} else { skip };
assert(l1 = l2)

#h1 = 0
 if (h1)
 #h1
$$\neq$$
 0
#l1=#l2 \land #h1 \neq 0,
[l1 \rightarrow #l1, h1 \rightarrow #h1, l2 \rightarrow #l2, h2 \rightarrow #h2]
#l1=#l2 \land #h1 \neq 0,
[l1 \rightarrow 1, h1 \rightarrow #h1, l2 \rightarrow #l2, h2 \rightarrow #h2]

$$|11 := 1$$

$$#l1=#l2 \land #h1 \neq 0,$$

$$[l1 \rightarrow 1, h1 \rightarrow #h1, l2 \rightarrow #l2, h2 \rightarrow #h2]$$

$$#h2 = 0$$

$$#h2 \neq 0$$

$$#l1=#l2 \land #h1 \neq 0 \land #h2 \neq 0$$

1

 $\#11=\#12 \land \#h1 \neq \emptyset \land \#h2 \neq \emptyset,$ [11>1, h1>#h1, 12>#12, h2>#h2]



#h2
$$\neq 0$$
 if (h2)
#h2 = 0
#l1=#l2 \wedge #h1 $\neq 0 \wedge$ #h2 = 0,
[l1 \rightarrow 1, h1 \rightarrow #h1, l2 \rightarrow #l2, h2 \rightarrow #h2]
skip
#l1=#l2 \wedge #h1 $\neq 0 \wedge$ #h2 = 0,
[l1 \rightarrow 1, h1 \rightarrow #h1, l2 \rightarrow #l2, h2 \rightarrow #h2]

Self-Composition – example 2 skip assume(11 = 12);if (h1) { $\#11=\#12 \land \#h1 \neq 0 \land \#h2 = 0,$ $\lceil 11 \rightarrow 1, h1 \rightarrow \#h1, 12 \rightarrow \#12, h2 \rightarrow \#h2 \rceil$ 11 := 1} else { skip }; assert(l1 = l2)if (h2) { $\#11=\#12 \land \#h1 \neq 0 \land \#h2 = 0 \land 1 \neq \#12$ **SAT**? 12 := 1} else { skip }; **YES!** [#l1→0, #h1→1, #l2→0, #h2→0] assert(11 = 12)





If we find a bug, we know that the program is **not** secure

But what if we do **not** find a bug?

Is the program secure?



But what if we do **not** find a bug?

The program is secure if we covered all the possible **execution paths**

How can we know that?



How do we know if we covered all possible execution paths?

The disjunction of all final path conditions must be True



The final symbolic store is always the same!

What is the SAT query?

[l1→**1+#y1**, l2→**1+#y2**, h1→#h1, h2→#h2, y1→#y1, y2→#y2, z1→#z1, z2→#z2] **SELF-COMPOSITION - FORMALLY**

Idea: Reduce non-interference to a **safety property** by transpiling the given program

$$s \in \mathcal{NI}(\Gamma) \iff \llbracket \mathcal{C}(s) \rrbracket \subseteq \mathcal{T}(\Gamma)$$

 $\mathcal{T}(\Gamma)$ - a **safety property** that only depends on Γ

C - a transpiler that computes the self-composition of s

Self-Composition – Formally

$$\frac{\operatorname{dom}(\theta_1) = \operatorname{dom}(\theta_2) = \operatorname{vars}(s) \quad \operatorname{rng}(\theta_1) \cap \operatorname{rng}(\theta_2) = \emptyset}{C(s) \triangleq \theta_1(s); \theta_2(s)}$$

$$\mathcal{T}(\Gamma) \triangleq \{ [\langle \rho_0, s_0 \rangle, ..., \langle \rho_n, s_n \rangle] \\ | s_n = \mathsf{skip} \land \forall_{0 \le i < n} \langle \rho_i, s_i \rangle \to \langle \rho_{i+1}, s_{i+1} \rangle \\ \land \big(\bigwedge \{ (\rho_0(\theta_1(x)) = \rho_0(\theta_2(x))) \mid \Gamma(x) = L \} \\ \Rightarrow \bigwedge \{ (\rho_n(\theta_1(x)) = \rho_n(\theta_2(x))) \mid \Gamma(x) = L \} \big) \}$$

SELF-COMPOSITION - FORMALLY

$$dom(\theta_1) = dom(\theta_2) = vars(s) \quad rng(\theta_1) \cap rng(\theta_2) = \emptyset$$

$$s_{assume} = assume \left(\bigwedge \{ (\theta_1(x) = \theta_2(x)) \mid \Gamma(x) = L \} \right)$$

$$s_{assert} = assert \left(\bigwedge \{ (\theta_1(x) = \theta_2(x)) \mid \Gamma(x) = L \} \right)$$

 $C(s) \triangleq s_{assume}; \theta_1(s); \theta_2(s); s_{assert}$

IV. CONCOLIC SYMBOLIC EXECUTION

Idea: Execute the program concretely **and** symbolically at the same time



P. Godefroid

Why?

- Symbolic execution is often too expensive...
- Back-end constraint solvers sometimes (often!) cannot find the answer: UNKNOWN

Concolic Testing: Main Algorithm

Input₀ = *Pick random vector* Coverage = False i = 0While (Input_i ≠ **NULL**) { (RES_i, PC_i)= Run Program with Input_i Coverage = Coverage $\vee PC_i$ $Input_{i+1} \in Models(\neg Coverage)$ i = i+1

z := 2*y; if (z = x) { if (x > y+10) { assert(false) } else { skip }; } else { skip }

Step 1:

Inputs₀ = [
$$x \rightarrow 22, y \rightarrow 7$$
]
(RES₁, PC₁) = (OK, ($x \neq 2*y$))
Coverage = ($x \neq 2*y$)
Inputs₁ = [$x \rightarrow 2, y \rightarrow 1$]

z := 2*y; if (z = x) { if (x > y+10) { assert(false) } else { skip }; } else { skip }

Step 2:

z := 2*y; if (z = x) { if (x > y+10) { assert(false) } else { skip }; } else { skip }

Step 3:

A LOT MORE TO COVER...

- 1. Symbolic execution with **data structures**
 - Lazy-Initialization
 - Data-structure **unfolding**
- 2. Declassification
- 3. Other Security Properties: Confinement

4. Invariants and verification

MASTER PROJECTS



Bounded model checking for TypeScript via symbolic execution and compilation

Code-stepping regular expressions in the browser

Building a symbolic execution engine for your favorite programming language

MASTER PROJECTS



Bounded model checking for JavaScript regular expressions

Symbolically debugging secure information flow in the browser

And **more...** Check my website: http://web.ist.utl.pt/**jose.fragoso**
MASTER PROJECTS



Potential collaboratons with:

Imperial College London



