Specification Inference for Explicit Information Flow Problems

# Merlin

Benjamin Livshits, Aditya V. Nori, Sriram K. Rajamani
Microsoft Research

Anindya Banerjee
IMDEA Software

# Mining Security Specifications

- **Problem:** Can we automatically infer which routines in a program are sources, sinks and sanitizers?

- **Technology:** Static analysis + Probabilistic inference

- **Applications:**
  – Lowers false errors from tools
  – Enables more complete flow checking

- **Results:**
  – Over 300 new vulnerabilities discovered in 10 deployed ASP.NET applications

Microsoft Research

# Motivation

# Static Analysis Tools for Security

- Web application vulnerabilities are a serious threat!

# Web Application Vulnerabilities



```
$username = $_REQUEST['username'];
$sql = "SELECT * FROM Students WHERE username = '$username';
```

Microsoft
**Research**

# Propagation graph

```
void ProcessRequest()
 {
    string s1 = ReadData1("name");
    string s2 = ReadData2("encoding");

    string s11 = Prop1(s1);
    string s22 = Prop2(s2);

    string s111 = Cleanse(s11);
    string s222 = Cleanse(s22);

    WriteData("Parameter " + s111);
    WriteData("Header " + s222);
 }
```



Propagation graph
m1→ m2 iff information flows "explicitly"
from m1 to m2

# Specification

# Vulnerability

- **Source**
  - returns tainted data
- **Sink**
  - error to pass tainted data
- **Sanitizer**
  - cleanse or untaint the input
- **Regular nodes**
  - propagate input to output

- *Every path from a source to a sink should go through a sanitizer*

- *Any source to sink path without a sanitizer is an information flow vulnerability*

Microsoft®
**Research**

# Information flow vulnerabilities

```
void ProcessRequest()
{
    string s1 = ReadData1("name");
    string s2 = ReadData2("encoding");

    string s11 = Prop1(s1);
    string s22 = Prop2(s2);

    string s111 = Cleanse(s11);
    string s222 = Cleanse(s22);

    WriteData("Parameter " + s111);
    WriteData("Header " + s222);
}
```

# Information flow vulnerabilities

# Information flow vulnerabilities

# Goal

*Given a propagation graph, can we infer a specification or 'complete' a partial specification?*

Assumption
*Most flow paths in the propagation graph are secure*

# Algorithms

# Merlin Architecture

# Propagation Graph Construction



```
void ProcessRequest()
{
    string s1 = ReadData1("name");
    string s2 =
ReadData2("encoding");

    string s11 = Prop1(s1);
    string s22 = Prop2(s2);

    string s111 = Cleanse(s11);
    string s222 = Cleanse(s22);

    WriteData("Parameter " + s111);
    WriteData("Header " + s222);
}
```

Microsoft **Research**

# Inference?

# Path constraints

- For every acyclic path $m_1\, m_2\, ...\, m_n$ the probability that $m_1$ is a source, $m_n$ is a sink, and $m_2, ..., m_{n-1}$ are not sanitizers is very low



Exponential number of path constraints: $O(2^{|V|})$!

# Triple constraints

- For every triple $\langle m_1, m_i, m_n \rangle$ such that $m_i$ is on a path from $m_1$ to $m_n$, the probability that $m_1$ is a source, $m_n$ is a sink, and $m_i$ is not a sanitizer is very low



Cubic number of triple constraints: $O(|V|^3)$!

# Minimizing Sanitizers

# Minimizing Sanitizers

For every pair of nodes $m_1$, $m_2$ such that $m_1$ and $m_2$ lie on the same path from a potential source to a potential sink, the probability that both $m_1$ and $m_2$ are sanitizers is low

# Need for probabilistic constraints



Triple constraints
- ¬(a ∧ ¬b ∧ d)
- ¬(a ∧ ¬c ∧ d)

Avoid double sanitizers
- ¬(b ∧ c)

- a ∧ d ⇒ b
- a ∧ d ⇒ c
- ¬(b ∧ c)

# Boolean formulas as probabilistic constraints

$$(x_1 \lor x_2) \land (x_1 \lor \neg x_3)$$

$$C_1 \qquad C_2$$

$$f(x_1, x_2, x_3) = f_{C1}(x_1, x_2) \land f_{C2}(x_1, x_3)$$

$$f_{C1}(x_1, x_2) = \begin{cases} 1 \text{ if } x_1 \lor x_2 = \text{true} \\ 0 \text{ otherwise} \end{cases}$$

$$f_{C2}(x_1, x_3) = \begin{cases} 1 \text{ if } x_1 \lor \neg x_3 = \text{true} \\ 0 \text{ otherwise} \end{cases}$$

# Boolean formulas as probabilistic constraints

$$(x_1 \lor x_2) \land (x_1 \lor \neg x_3)$$

$$\underbrace{\phantom{(x_1 \lor x_2)}}_{C_1} \quad \underbrace{\phantom{(x_1 \lor \neg x_3)}}_{C_2}$$

$$f(x_1, x_2, x_3) = f_{C1}(x_1, x_2) \land f_{C2}(x_1, x_3)$$

$$f_{C1}(x_1, x_2) = \begin{cases} 1 \text{ if } x_1 \lor x_2 = \text{true} & 0.9 \\ 0 \text{ otherwise} & 0.1 \end{cases}$$

$$f_{C2}(x_1, x_3) = \begin{cases} 1 \text{ if } x_1 \lor \neg x_3 = \text{true} & 0.9 \\ 0 \text{ otherwise} & 0.1 \end{cases}$$

$$p(x_1, x_2, x_3) = f_{C1}(x_1, x_2) \times f_{C2}(x_1, x_3) / Z$$

$$Z = \sum_{x1, x2, x3} (f_{C1}(x_1, x_2) \times f_{C2}(x_1, x_3))$$

# Solution = Marginalization

marginal

$$p_i(x_i) = \sum_{x1} \dots \sum_{x(i-1)} \sum_{x(i+1)} \dots \sum_{xN} p(x_1, \dots, x_N)$$

- Step 1: choose $x_i$ with highest $p_i(x_i)$ and set $x_i$ = true if $p_i(x_i)$ is greater than a threshold, false otherwise
- Step 2: recompute marginals and repeat Step 1 until all variables have been assigned

# Factor graphs: efficient computation of marginals



$$f_{C1}(x_1, x_2) = \begin{cases} 1 \text{ if } x_1 \vee x_2 = \text{true} \\ 0 \text{ otherwise} \end{cases}$$

$$f_{C2}(x_1, x_3) = \begin{cases} 1 \text{ if } x_1 \vee \neg x_3 = \text{true} \\ 0 \text{ otherwise} \end{cases}$$

# Factor Graphs

# Probabilistic Inference

|  | Source | Sanitizer | Sink |
|---|---|---|---|
| ReadData1 | .95 | .001 | .001 |
| ReadData2 | .5 | .5 | .5 |
| Cleanse | .5 | .5 | .5 |
| WriteData | .5 | .5 | .85 |
| ... |  |  |  |

Probabilistic Inference

|  | Source | Sanitizer | Sink |
|---|---|---|---|
| ReadData1 | .95 | .001 | .001 |
| ReadData2 | .5 | .5 | .5 |
| Cleanse | .01 | .997 | .03 |
| WriteData | .5 | .5 | .85 |
| ... |  |  |  |

# Paths vs. Triples

Path$(G = \langle V, E \rangle)$
Returns:
Mapping $m$ from $V$ to the set $\{0, 1\}$

1: for all paths $p = s, \ldots, n$ from potential sources to sinks in $G$ do
2:    assume( $m(p) \notin 10^*1$ ) $\oplus_{c_p}$ assume($m(p) \in 10^*1$)
3: end for
Post expectation: $[\forall$ paths $p$ in $G, m(p) \notin 10^*1]$.

Theorem
Path refines Triple

Triple$(G = \langle V, E \rangle)$
Returns:
Mapping $m$ from $V$ to the set $\{0, 1\}$

1: for all triples $t = \langle s, w, n \rangle$ such that $s$ is a potential source, $n$ is a potential sink
     and $w$ lies on some path from $s$ to $n$ in $G$ do
2:    assume( $m(\langle s, w, n \rangle) \neq 101$) $\oplus_{c_t}$ assume( $m(\langle s, w, n \rangle) = 101$)
3: end for
Post expectation: $[\forall$ paths $p$ in $G, m(p) \notin 10^*1]$.

# Experiments

# Implementation

- Merlin is implemented in C#
  - Uses CAT.NET for building the propagation graph
  - Uses Infer.NET for probabilistic inference
    - http://research.microsoft.com/infernet

# Experiments

10 line-of-business applications written in C# using ASP.NET



| Type | Count | Revisions |
|------|-------|-----------|
| Sources | 27 | 16 |
| Sinks | 77 | 8 |
| Sanitizers | 7 | 2 |

# Summary of Discovered Specifications

# Summary of Discovered Vulnerabilities

# Experiments - summary

- 10 large Web apps in .NET
- Time taken per app < 4 minutes
- New specs: 167
- New vulnerabilities: 322
- False positives removed: 13
- Final false positive rate for CAT.NET after Merlin < 1%

# Summary

- **Merlin** is first practical approach to infer explicit information flow specifications
- Design based on a formal characterization of an approximate probabilistic constraint system
- Able to successfully and efficiently infer explicit information flow specifications in large applications which result in detection of new vulnerabilities

[http://research.microsoft.com/merlin](http://research.microsoft.com/merlin)

Microsoft®
**Research**