# **Static Taint Analysis for C with LLVM**

Xavier Noumbissi

Department of Electrical and Computer Engineering
University of Waterloo

## Problem: Prevent Software Vulnerabilities

- Format String Attacks

- SQL Injection

- Cross Site Scripting, etc.

# Solution: Track Used of Non Trusted Program Input

- Non trusted program input: <span style="color:red">Tainted Input</span>

- Source: origin of tainted input

- Sink: use of tainted input

- Taint Propagation: tracking tainted input

# Example

```
1   int main() {
2     int x, b1, b2, y;
3     scanf("%d", &x);
4     b1 = even(x);
5     b2 = odd(3);
6     y = compute(x);
7     return 0;
8   }

10  int compute(int x) {
11    int sum, i;
12    if (x == 2)
13      scanf("%d", &sum);
14    else
15      sum = 0;
16    for(i = 0; i < x; ++ i)
17      sum += i;
18    return sum;
19  }

21  int odd(int x) {
22    if (x == 1)
23      return 0;
24    else
25      return even(x - 1);
26  }

28  int even(int x) {
29    if (x == 0)
30      return 1;
31    else
32      return odd(x - 1);
33  }
```

Figure 1. Motivating Example

# Solution: In This Project

- Implicit Taint Propagation: due to Control Flow

- Explicit Taint Propagation: due to Data Flow

**Expected Contributions**

- Algorithm to statically detect use of tainted values in C programs

- Handling of interprocedural taint propagation

- Implementation of the algorithm in LLVM

# Taint Analysis: Relevant C Program Statements

| Statement Type | C Code |
| --- | --- |
| COPY | $p = q$ |
| LOAD | $p = *q$ |
| STORE | $*p = q$ |
| ADDROF | $p = \&a$ |
| SOURCE | *call gets* |

# Taint Analysis: Transfer Functions

- COPY [$p = q$]: taint p iff q is tainted

- LOAD [$p = *q$]: taint p iff it exists $t_q = *q \wedge t_q$ is tainted

- STORE [$*p = q$]: nothing

- ADDROF [$p = \&a$]: nothing

- SOURCE [*call gets(p)*]: taint all $t_q$ s.t $t_q = *p$

# Interprocedural, Context-Sensitive Analysis

- Analysis of a callee start with the taint
  assumptions from the caller

# Taint Information from Source

- Developer specify sources and sinks in configuration file

- Analysis do not analyze sources and sinks

- Analysis use annotations for sources: taint propagation

# Algorithm Analyze: Implements Interprocedural Analysis

```
input  : func : Proc,
         initDataFlow : Inst → (Var → 2^Inst)
output:
1  s_0 ← first(f)
2  input[s_0] ← initDataFlow(s_0)
3  worklist ← {s_0}
4  while worklist ≠ ∅ do
5      i ← next(worklist)
6      output[i] ← Flow(Analyze, i)
7      foreach j ∈ succs(i) do
8          if output[i] ≠ input[j] then
9              input[j] ← input[j] ⊔ output[i]
10             worklist ← worklist ∪ {j}
11         end
12     end
13 end
```

Algorithm 1: Analyze

# Algorithm Flow: Implements Transfer Functions

```
    input : caller : Proc, s : Inst
    output:
1  switch TypeOf(s) do
2  |  case COPY [p = q]
3  |  |   if inFlowₛ[q] = ∅ then
4  |  |   |   outFlowₛ[p] ← inFlowₛ[q] ∪ {s};
5  |  |   end
6  |  endsw
7  |  case LOAD [p = *q]
8  |  |   foreach a ∈ pt_{[s]}(q) do
9  |  |   |   if inFlowₛ[a] = ∅ then
10 |  |   |   |   outFlowₛ[p] ← inFlowₛ[p] ∪ {s}
11 |  |   |   end
12 |  |   end
13 |  endsw
14 |  case SOURCE [call func(a₀, a₁, ..., aₙ)]
15 |  |   foreach k ∈ {0, 1, ..., n} do
16 |  |   |   if taint(k) then
17 |  |   |   |   outFlowₛ[aₖ] ← inFlowₛ[aₖ] ∪ {s}
18 |  |   |   end
19 |  |   end
20 |  endsw
21 |  case CALL [call func(a₀, a₁, ..., aₙ)]
22 |  |   if caller = func then
23 |  |   |   foreach aₖ, k ∈ {0, 1, ..., n} do
24 |  |   |   |   fₖ ← formal (func, aₖ)
25 |  |   |   |   if aₖ ∈ P then
26 |  |   |   |   |   foreach b ∈ pt_{[s]}(aₖ) do
27 |  |   |   |   |   |   inFlowₛ[fₖ] ← inFlowₛ[fₖ] ∪ inFlowₛ[b]
28 |  |   |   |   |   end
29 |  |   |   |   end
30 |  |   |   |   else if aₖ ∈ A then
31 |  |   |   |   |   tₓ ← toplevel (aₖ)
32 |  |   |   |   |   foreach b ∈ pt_{[s]}(tₓ) do
33 |  |   |   |   |   |   inFlowₛ[fₖ] ← inFlowₛ[fₖ] ∪ inFlowₛ[b]
34 |  |   |   |   |   end
35 |  |   |   |   end
36 |  |   |   end
37 |  |   |   Flow(caller, func)
38 |  |   end
39 |  endsw
40 |  case ADDROF [p = &a]
41 |  case STORE [*p = q]
42 |  case SINK [call func]
43 |  endsw
```

## Implementation

- LLVM infrastructure is ready

- Need to implement the analysis

## Future Work

- Make analysis modular

- Better handling of (mutual-) recursive function calls

# Conclusion

- Implementation should take less than 1 months

- Need to do evaluation on real world C programs

- We are optimistic about future results

**Thank You!**

# Comments & Questions