

Visualizing Information Flow through C Programs

Automated Security Analysis Project

Joe Hurd, Aaron Tomb and David Burke

Galois, Inc.

{joe,atomb,davidb}@galois.com

HCSS

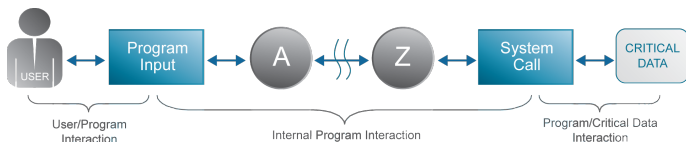
May 10, 2010

Talk Plan

- 1 Introduction
- 2 Information Flow Analysis
- 3 C Information Flow Tool (Cift)
- 4 Summary

Software Security Evaluation

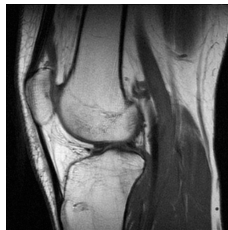
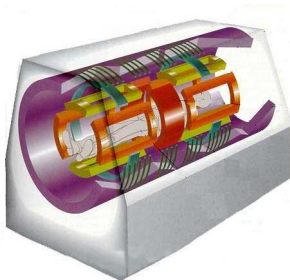
- Evaluating the security of a program generally focuses on:
 - The **attack surface** (e.g., the interface to the user or network).
 - The **critical data** (e.g., crypto keys, database queries).
- **Typical Question:** What are the possible effects of changes at the attack surface on the critical data?



- Answering this requires an understanding of how information flows through the program.

Information Flow Diagnostic Tool

- Tight time constraints mean that evaluators often cannot look at every line of the codebase.
- **Project Goal:** Develop an interactive diagnostic tool that allows an evaluator to scan for anomalies in the program information flows.



Information Flow and Security Properties

- Many program security properties can be expressed in terms of potential information flow between program variables.
- **Confidentiality:** ✓
 - There are no information flows from **secret** variables to **public** variables.
- **Integrity:** ✓
 - There are no information flows from **tainted** variables to **critical** variables.
- **Availability:** ✗
 - No way to specify that a flow **must** happen.

Information Flow and Confidentiality

- **Use Case:** Ensure Bell-La Padula properties hold for a cross domain application.
- **Annotate:** Program variables with their sensitivity level.
- **Check:** There are no flows from higher sensitivity to lower sensitivity variables.
- **Example:**

Code (Confidentiality Bug)

```
void f() {  
    int k = get_secret_key();  
  
    publish_to_internet(k);  
}
```

Information Flow and Integrity

- **Use Case:** Defend against SQL injection attacks.
- **Annotate:** Tainted data variables; critical data variables; and validation functions.
- **Check:** All flows from tainted variables to critical variables go through validation functions.
- **Example:**

Code (Integrity Check Succeeds)

```
void f() {  
    int data = get_user_input();  
  
    data = validate_input(data);  
  
    query_sql_database(data);  
}
```

User Centered Design

- “*All tools are user interfaces*” – Clark Dodsworth
- **Evaluator/Tool Workflow:**
 - ① The evaluator seeds the analysis by annotating some program variables as sensitive data or dangerous user input.
 - ② The tool uses the annotations to find candidate insecure information flows.
 - ③ The evaluator examines the flows, and removes false positives by providing additional annotations so that the tool can make a more precise analysis.
- **Tool Requirements:**
 - ① Scalable analysis of program information flow.
 - ② Intuitive visualization of information flow in terms of source code.

Analysis Evidence

- **Evidence of Security Bugs:** Insecure information flows are presented as a sequence of assignments on the control flow.
 - **False Positives:** The evaluator uses the tool to browse the insecure information flows, and adds annotations to eliminate false positives.
- **Evidence of Assurance:** The analysis computes a conservative over-approximation of information flow on a subset of the programming language.
 - **False Negatives:** The tool will emit a warning message when the analysis detects that the program is outside of the conservative subset, allowing the evaluator to assess the residual risk.

Static Analysis

- Static analysis is a program verification technique that is complementary to testing.
 - Testing works by **executing** the program and checking its run-time behavior.
 - Static analysis works by examining the **text** of the program.
- Driven by new techniques, static analysis tools have recently made great improvements in scope.
 - **Example 1:** Modern type systems can check **data integrity** properties of programs at compile time.
 - **Example 2:** Abstract interpretation techniques can find memory problems such as **buffer overflows** or **dangling pointers**.
 - **Example 3:** The TERMINATOR tool developed by Microsoft Research can find **infinite loops** in Windows device drivers that would cause the OS to hang.

Information Flow Static Analysis: Requirements

- **Evidence:** Generating evidence of assurance relies on the information flow static analysis being **sound**:
 - 1 Define a sound static analysis on a simple input language.
 - 2 Implement a conservative translator from the target programming language to the simple input language.
- **Scalability:** To help the static analysis scale up to realistically sized codebases, we design it to be **compositional**.
 - Preserve function calls in the input language.
- **Program Understanding:** The analysis result must help an evaluator understand how information flows through the program **source code**.
 - Link each step in the analysis to the program source code.

The Input Language

- **Variables**

- Global variables.
- Local variables of a function.

- **Statements**

- Simple variable assignment $v_1 \leftarrow v_2$.
- Function call $v \leftarrow f(v_1, \dots, v_n)$.

- **Functions**

- Special local variables representing input arguments $\$arg1, \dots, \$argN$ and return value $\$ret$.
- A function contains a **set** of statements (flow insensitive).

- **Programs**

- A set of functions, including a distinguished **main** function where execution begins.

The Input Language as a C Subset

Code (Example Program)

```
/* High global variables */
int high_in; int high_out;

/* Low global variables */
int low_in; int low_out;

int f(int x) { return x + 1; }

int main() {
    high_in = 42;
    low_in = 35;

    high_out = f(high_in);
    low_out = f(low_in);

    return 0;
}
```

Program Information Flow Graph

- The **program information flow graph** G is a directed graph between global variables.
 - An edge $x \rightarrow y$ indicates that the program enables a possible information flow from x to y .
 - Can be efficiently computed using abstract interpretation.
- **Key Property:** The information flow analysis is **sound**.
- **Key Property:** The analysis is **context sensitive**.

Analysis (Example Program Information Flow Graph)

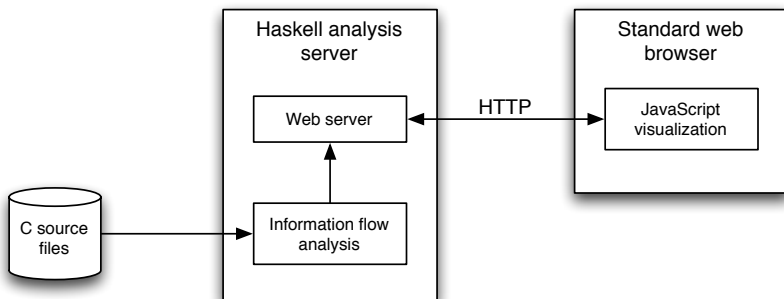
$$G = \{\text{low_in} \rightarrow \text{low_out}, \\ \text{high_in} \rightarrow \text{high_out}\}$$

Information Flow Static Analysis of C Code

- The following front end processing is performed to translate C code to the input language:
 - ① **Preprocessing:** The C preprocessor.
 - ② **Parsing:** The Haskell Language.C package.
 - ③ **Simplification:** Normalizing expressions (like CIL).
 - ④ **Variable Classification:** Special handling for address-taken locals and dynamically allocated memory.
 - ⑤ **Pointer Analysis:** Anderson's algorithm replaces each indirect reference with a set of direct references.
- **Key Property:** The front end processing is **conservative**.
 - Every information flow in the C code is translated to an information flow in the input language.
 - **Assumption:** the C code is memory safe.

Cift Architecture

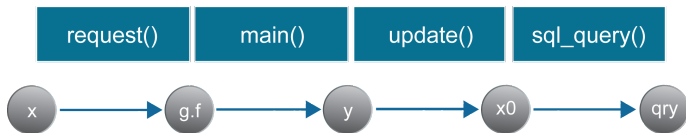
The [C Information Flow Tool](#) (Cift) allows evaluators to examine information flows in C code using a [standard web browser](#).



The architecture is designed to support [multiple simultaneous users](#) browsing code and [sharing annotations](#).

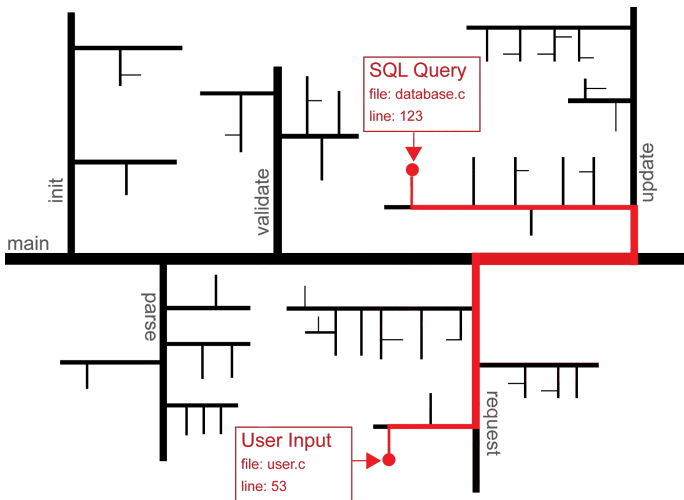
Visualizing Information Flow

- A program information flow consists of many assignments distributed across the codebase:



- Tracking a long information flow across source code involves much tedious opening, closing and searching of files.
 - *“Evaluating software is like frying 1,000 eggs”*
- A different visualization solution is needed.

Right-Angle Fractal Call Trees

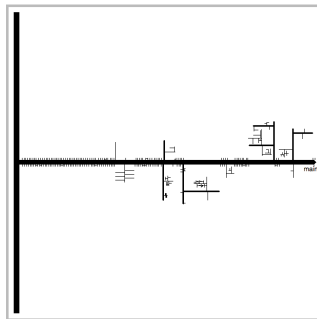
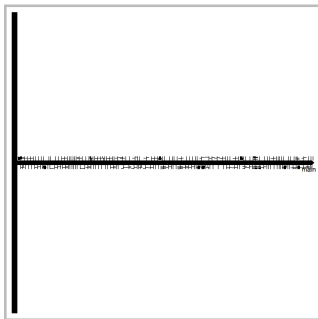


Restricting to Variables of Interest

- **Problem:** A typical large program contains many variables V in the program information flow graph: [information overload](#).
- **Solution:** Allow the user to specify a subset $X \subseteq V$ of interesting variables.
- Remove the uninteresting global variables $V - X$ from the program information flow graph one by one.
 - When removing a variable v , an extra edge $x \rightarrow y$ must be added between every pair of variables x, y satisfying $x \rightarrow v$ and $v \rightarrow y$.
 - This amounts to computing the transitive closure of the program information flow graph on demand.
- A first step towards information flow annotations.

Emphasizing Call Tree Paths of Interest

Problem: Functions with too many function calls result in an uninformative **hairy spike** (left graphic).



Solution: Emphasize function calls contributing to information flows between **variables of interest** (right graphic).

Open Source Benchmarks

- All experiments were carried out on a MacBook Pro 2.2Ghz Core 2 Duo with 4Gb of RAM, using GHC 6.12.1.
- Analyzing the 67 KLoC C implementation of OpenSSH takes 1:53s of CPU time and consumes 1.6Gb of RAM.
- Analyzing the 94 KLoC C implementation of the SpiderMonkey JavaScript interpreter takes 6:49s of CPU time and consumes 1.3Gb of RAM.

Cift Development Plan

Milestones:

- ✓ Develop an automatic information flow analysis that scales up to realistic C codebases.
- ✓ Develop a visualization technique for program information flow that is grounded in the source code.
- ✓ Implement a research prototype tool to examine information flows in C programs.
- Develop an annotation language for information flow properties of C functions and variables.
- Allow users to edit annotations through the browser interface and see the resulting effects on the analysis.

Future Plans

- Extend the scope of the information flow analysis.
 - Supporting array sensitivity to distinguish the elements of an array or cells in a memory block.
 - Adding flow sensitivity and a clobber analysis to detect failures to sanitize confidential data after use.
 - Target LLVM to extend the analysis to C++/Ada/etc.
- Support higher-level information flow specifications.
 - Derive program specifications from higher-level security policies.
 - Track information flow across module and language barriers.

Summary

- **This Talk:** We have presented a research prototype static analysis tool that an evaluator can use to visualize how information flows through C programs.
- **Feedback Welcome:** Would a tool like this make your job easier? Please let us know what features you'd like to see in a program understanding tool.

joe@galois.com

