

Interactively Verifying Absence of Explicit Information Flows in Android Apps

Abstract

App stores are increasingly the preferred mechanism for distributing software, including mobile apps (Google Play), desktop apps (Mac App Store and Ubuntu Software Center), computer games (the Steam Store), and browser extensions (Chrome Web Store). The centralized nature of these stores has important implications for security. While platforms have unprecedented ability to audit apps, users now trust hosted apps, making them more vulnerable to malware that evades detection and finds its way onto the app store. Sound static explicit information flow analysis has the potential to significantly aid the auditor, but it is handicapped by high false positive rates. Instead, auditors currently rely on a combination of dynamic analysis (which is unsound) and lightweight static analysis (which cannot identify information flows) to help detect malicious behaviors.

We propose a process for producing apps certified to be free of malicious explicit information flows. In practice, imprecision in the reachability analysis is a major source of false positive information flows that are difficult to understand and discharge. In our approach, the developer provides tests that specify what code is reachable, allowing the static analysis to restrict its search to tested code. The app hosted on the store is instrumented to enforce the provided specification (i.e., executing untested code terminates the app). We use abductive inference to minimize the necessary instrumentation, and then interact with the developer to ensure that the instrumentation only cuts unreachable code. We demonstrate the effectiveness of our approach in verifying a corpus of 78 Android apps—our interactive verification process successfully discharges 11 out of the 12 false positives.

Categories and Subject Descriptors F.3.2 [*Semantics of Programming Languages*]: Program analysis; F.3.1 [*Specifying and Verifying and Reasoning about Programs*]: Mechanical verification

Keywords interactive verification; abductive inference; specifications from tests

1. Introduction

Android malware has become increasingly problematic as the popularity of the platform has skyrocketed in the past few years [55]. App stores currently identify malware using a two-step process: first, they use an automated *malware detection pipeline* to flag suspicious apps, and then a human auditor manually reviews flagged apps. The detection pipeline typically combines dynamic analysis (e.g., dynamic information flow [16]) and static analysis (e.g., identifying execution of dynamically loaded code [21]); subsequent manual examination is necessary both because the static analysis may be imprecise and because determining if suspicious behavior is malicious requires understanding the behavior’s purpose (e.g., a map app should not send location data to the Internet, but it is okay for a location sharing app to do so). For example, Google uses a combination of dynamic analysis [7] and manual analysis [6] to audit Android apps on Google Play.

Important families of malware include apps that leak location, device ID, or contact information, and malware that covertly send SMS messages to premium numbers [18, 54]. Static and dynamic analyses for finding *explicit information flows* (i.e., flows arising only from data dependencies [41]) can be used to identify such malware [8, 16, 18, 19]. While explicit information flows do not describe all malicious behaviors, we focus on them because they can be used to characterize a significant portion (and possibly the majority) of existing malware (especially in conjunction with lightweight static analyses such as [21]).

The drawback of dynamic explicit information flow analysis is the potential for false negatives. Malware can easily avoid detection by making the relevant information flow difficult to trigger. For example, the malware can detect if the execution environment is an emulator, and if so disable malicious behaviors [49]. While static explicit information flow analysis avoids this problem, it can have a high false positive rate. Because of their global nature, false positive information flows can be time-consuming for a human auditor to understand and discharge. Furthermore, the rate at which the auditor can evaluate suspicious applications is likely the bottleneck for finding malware. This means that very few benign apps should be flagged for manual review—i.e., the malware detection pipeline must have a low false positive rate. As a consequence, app stores currently rely on approaches that achieve low false positive rates (possibly at

the expense of higher false negative rates), thereby excluding static explicit information flow analysis.

Our goal is to design a sound verification process for enforcing the absence of malicious explicit information flows that guarantees no false negatives and maintains a low false positive rate. Our first key observation is that for static explicit information flow analysis, a large number of false positives are flows through unreachable code. Such false positives result from imprecision in the static reachability analysis, which can be a major problem for static analyses [9, 14]. In our setting, this imprecision is caused both by an imprecise callgraph (due to virtual method calls) and by the lack of path sensitivity (for scalability reasons). Furthermore, manually written specifications are needed to describe the entry points of an Android app, and missing specifications become a huge source of false negatives. Oftentimes, the unreachable code is found in large third-party libraries used by the app. In our experiments, 92% of false positives were flows through unreachable code.

Our second key observation is that currently the burden of identifying and discharging false positives is entirely placed on the auditor, despite the fact that the developer is most familiar with the app’s code. Our approach shifts some of this burden onto the developer: we require that the developer specify which methods are reachable in the app. These reachability specifications allow the static analysis to restrict its search space to reachable code, thereby reducing the false positive rate.

In practice, we envision that the developer will provide reachability specifications by supplying tests that exercise the app code—the specification we extract is that only tested code is reachable. Using tests as specifications has a number of advantages. First, developers routinely write tests, so this approach both leverages existing tests and gives developers a familiar interface to the specification process. Second, tests are executable, which means that the auditor can verify the correctness of the specifications (otherwise, the developer can specify all code as reachable, eliminating the benefits of our approach). Finally, concrete test cases can benefit the auditor in case the app must be manually examined. For our technical development, we assume that specifications are extracted from tests, though any method for obtaining correct reachability specifications suffices.

Of course, a malware developer can attempt to evade detection by specifying that nothing is reachable. Our solution is simple: we enforce the developer-provided specifications by instrumenting the app to terminate if code not specified to be reachable (e.g., not covered by any of the developer-provided tests) is actually reached during runtime. The instrumented app is both consistent with the developer’s specifications, and statically verified to be free of explicit information flows.

In practice, it may be difficult for developers to provide tests covering all reachable code. Therefore, we take an it-

erative approach to obtaining tests. To enforce the security policy, it is only necessary to terminate the app if it reaches uncovered code that may also lead to a malicious explicit information flow. Rather than instrument all untested statements, we find a minimum size set of untested statements (called a *cut*) such that instrumenting these statements to terminate execution produces an app that is free of explicit information flows, and then propose this cut to the developer. If the developer finds the cut unsatisfactory, then the developer can provide new tests (or other reachability information) and repeat the process; otherwise, if the developer finds the cut satisfactory, then the cut is enforced via instrumentation as before. This process repeats until either a satisfactory cut is produced, or no valid cut exists (in which case the auditor must manually review the app). We call this process *interactive verification*.

If the developer allows (accidentally or maliciously) reachable code to be instrumented, then it may be possible for the app to terminate during a benign execution. To make the process more robust against such failures, we can produce *multiple, disjoint* cuts. We then instrument the program to terminate only if at least one statement from every cut is reached during an execution of the app.

More formally, our goal is to infer a predicate λ (the statements in the cut) such that the security policy ϕ (lack of explicit information flows) holds provided λ holds. We compute λ using *abductive inference* [15]: given facts χ (extracted from the app P) and security policy ϕ , abductive inference computes a minimum size predicate λ such that (i) $\chi \wedge \lambda \models \phi$ (i.e., λ together with the known program facts χ suffice to prove the security policy ϕ) and (ii) $\text{SAT}(\chi \wedge \lambda)$ (i.e., λ is consistent with the known program facts χ). In our setting, we augment χ with facts extracted from the tests. Then the app P is instrumented to ensure that λ holds, producing a verified app P' . Finally, we extend this process to infer multiple disjoint cuts $\lambda_1, \dots, \lambda_n$, and instrument P to terminate only if every λ_i fails.

We propose a novel algorithm for solving the abductive inference problem for properties formulated in terms of context-free language (CFL) reachability. The security policy ϕ states that certain vertices in a graph representation of the program are unreachable. Our key insight is to formulate the CFL reachability problem as a constraint system, which we encode as an integer linear program (ILP). Finding minimum cuts in turn corresponds to a minimum solution for the ILP. Our work has three main contributions:

- We formalize the notion of *interactive verification* for producing verified apps using abductive inference (Section 3).
- For properties ϕ formulated in terms of CFL reachability (Section 4), we reduce the abductive inference problem to an integer linear program (Section 5).

```

1. void leak(boolean flag, String data) {
2.     // @Sink("sendHTTP.param")
3.     if (flag) sendHTTP(data); }
4. // @Entry("onCreate")
5. void onCreate() {
6.     // @Source("getLocation.return")
7.     String loc = getLocation();
8.     Runnable runMalice = new Runnable() {
9.         void run() { leak(true, loc); }}
10.    Runnable runBenign = new Runnable() {
11.        void run() { leak(false, loc); }}
12.    runBenign.run(); }

```

Figure 1. A potential false positive information flow through program P_{onCreate} .

- We implement our framework (Section 6) for producing Android apps verified to be free of explicit information flows, and show that our approach scales to large Android apps, some with hundreds of thousands of lines of bytecode (Section 7).

2. Motivating Example

Consider the Java-like code shown in Figure 1, which we call P_{onCreate} . The goal is to prove $\phi_{\text{flow}} = \nexists(\text{source flows to sink})$. We assume that the information sources and sinks are given (or inferred, see [28]); in this example, the user’s location is the source, and the Internet is the sink. As discussed earlier, one major source of false positive information flows is unreachable code. We therefore improve the precision of our analysis by removing parts of the program that are statically proven to be unreachable.

However, statically computed reachability information can be very imprecise. One source of imprecision is from imprecision in the static callgraph. For example, using a callgraph generated using class hierarchy analysis, the analysis will not determine that line 9 cannot call `runMalice.run`. Even with a more precise callgraph, the static analysis may not be able to prove that `flag` is false in every execution. Hence a sound static analysis fails prove that ϕ_{flow} holds.

Our approach is to search for a subset X of program statements such that if every statement in X is unreachable, then ϕ_{flow} holds. Let S be the set of program statements, and let $S^* \subseteq S$ be the set of reachable program statements. The static analysis produces an approximate set of reachable statements \hat{S} satisfying $S^* \subseteq \hat{S} \subseteq S$. We infer predicates of the form $\lambda_X = \bigwedge_{s \in X} (s \notin S^*)$, where $X \subseteq \hat{S}$. In other words, λ_X specifies that the statements in X are unreachable. Any of the following choices for X would allow the static analysis to prove ϕ_{flow} for P_{onCreate} : (i) $\{3.\text{sendHTTP}(\text{data})\}$, (ii) $\{7.\text{getLocation}()\}$, and (iii) $\{9.\text{leak}(\text{true}, \text{loc})\}$.

Not all of these choices are desirable. By executing `onCreate`, we observe that `7.getLocation()` is reachable, so our inference algorithm returns either (i) or (iii). Suppose

(i) is returned; then our algorithm instruments P_{onCreate} to enforce that `3.sendHTTP(data)` is unreachable, producing a program P'_{onCreate} for which ϕ_{flow} provably holds. Because $(3.\text{sendHTTP}(\text{data}) \notin S^*)$ happens to be true for P_{onCreate} , we know P'_{onCreate} is semantically equivalent to P_{onCreate} . Furthermore, the instrumentation in P'_{onCreate} incurs no runtime overhead, since it is in fact unreachable. In general, we cannot guarantee that our method will always cut false positive flows only in unreachable code, but our experiments show that such an outcome is likely, especially when using multiple cuts (see Section 7).

2.1 Analyzing Callbacks

In addition to imprecision in the callgraph, another source of imprecision is whether to treat `runMalice.run` as a *callback*. Much of an Android app’s functionality is executed via callbacks that are triggered when certain system events occur, so callbacks must be annotated as program entry points. The Android framework provides thousands of callbacks, which can be specified either programmatically or in an XML file known as the Android manifest. Furthermore, many of these callbacks are undocumented, which makes manually annotating callbacks a time consuming and error prone task. If a callback annotation is missing, then reachable code may be excluded from the analysis, introducing unsoundness.

On the other hand, every callback must override an Android framework method—we call any such method a *potential callback*. Of course, not every potential callback is a true callback, for example any method overriding `Object.equals` is a potential callback but not a true callback. In our analysis, we make the sound assumption that *every* potential callback is a callback—that is, we conservatively overestimate the set of callbacks. We then infer a cut λ_X as before. For example, in Figure 1, the static analysis treats `runMalice.run` as a potential callback, and thus reports the flow of location data to the Internet. The abductive inference algorithm can return the cut $\lambda_X = (3.\text{sendHTTP}(\text{data}) \notin S^*)$, which as before guarantees that the program is free of explicit information flows.

While more precise analyses such as [13] exist for soundly identifying callbacks, they are still overapproximations, and furthermore may be prone to false negatives (e.g., failing to handle native code). Our approach is both simple to implement and sound.

3. Interactive Verification

We model the imprecision of static analysis by categorizing program facts as *may-facts* and *must-facts*. May-facts are facts that the static analysis cannot prove are false for all executions. For example, $(3.\text{sendHTTP}(\text{data}) \in S^*)$ is a may-fact for P_{onCreate} . Conversely, must-facts χ are facts that are shown to hold for at least one concrete execution of P . For example, since `7.getLocation()` is executed

by running `onCreate`, this is a must-fact for P_{onCreate} , i.e. $\chi \Rightarrow (7.\text{getLocation}() \in S^*)$.

Our static analysis takes as input a cut λ that asserts that some may-facts are false; for example, $\lambda_{9,11} = (9.\text{leak}(\text{true}, \text{loc}) \notin S^*) \wedge (11.\text{leak}(\text{false}, \text{loc}) \notin S^*)$ is a cut with which the static analysis can verify ϕ_{flow} for P_{onCreate} . These assumptions have a (finite) lattice structure $(\Lambda, \leq, \top, \perp)$, where $\lambda \leq \mu$ means: if $\chi \wedge \lambda \models \phi$ holds, then $\chi \wedge \mu \models \phi$ holds as well (i.e., μ makes stronger assumptions than λ). For example, $(9.\text{leak}(\text{true}, \text{loc}) \notin S^*) \leq \lambda_{9,11}$ because $\lambda_{9,11}$ makes stronger assumptions. The predicate \perp corresponds to no assumptions (and is guaranteed to hold), and \top corresponds to assuming that all may-facts are false.

Given $\lambda \in \Lambda$, our static analysis tries to prove $\chi \wedge \lambda \models \phi$ (i.e., it tries to prove ϕ assuming λ). When can we hope to find a valid cut λ that helps the static analysis prove ϕ ? Consider three cases:

1. The static analysis proves $\chi \wedge \perp \models \phi$. Since \perp always holds, the static analysis has proven that ϕ holds.
2. The static analysis cannot prove $\chi \wedge \perp \models \phi$, but proves $\chi \wedge \top \models \phi$. In this case, we can search for a valid cut $\lambda \in \Lambda$ with which the static analysis can prove ϕ .
3. The static analysis cannot prove $\chi \wedge \top \models \phi$. This means that even making best-case assumptions, the static analysis fails to prove ϕ , so no input $\lambda \in \Lambda$ can help prove ϕ .

In the first case, the app is already free of malicious information flows. In the third case, the app must be sent to the auditor for manual analysis. The second case is our case of interest, which we describe in more detail in the subsequent sections.

3.1 Abductive Inference

Our goal is to find a valid cut $\lambda \in \Lambda$ with which the static analysis can verify that the policy ϕ holds. Our central tool will be a variant of abductive inference where the known-facts χ are extracted from tests:

DEFINITION 3.1. Given must-facts χ extracted from dynamic executions, the *abductive inference problem* is to find a cut $\lambda \in \Lambda$ such that

$$\chi \wedge \lambda \models \phi \text{ and } \text{SAT}(\chi \wedge \lambda). \quad (1)$$

Additionally, we constrain λ to be *minimal*, i.e. there does not exist $\mu \in \Lambda$ satisfying (1) such that $\mu < \lambda$.

Abductive inference essentially allows us to compute minimal specifications λ that are simultaneously consistent with the must-facts χ and verify the policy ϕ .

We assume access to an oracle who we can query to obtain the tests used to extract must-facts χ :

DEFINITION 3.2. An *oracle* \mathcal{O} is a function that, on input specifications λ and program P , returns a test T_{new} showing that λ does not hold for P , or returns \emptyset if λ holds for P .

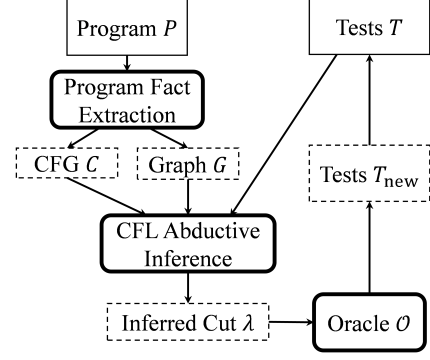


Figure 2. The interactive verification system.

In our setting, the oracle is the developer. If the developer is not satisfied with an inferred cut λ , then the developer can produce a new test case such that the extracted must-facts χ_{new} demonstrate that λ is an invalid specification for P , i.e. $\text{UNSAT}(\chi_{\text{new}} \wedge \lambda)$. The tool will update the must-facts $\chi \leftarrow \chi \wedge \chi_{\text{new}}$ and rerun the inference procedure. This process repeats until the developer is satisfied with λ , upon which verification is complete. This process is performed by the refinement loop in function `INTERACTIVECUT(P, T)` in Algorithm 1.

3.2 Instrumenting Cuts

Given cut λ , our framework produces $P' \leftarrow \text{INSTRUMENT}(P, \lambda)$, where P' is instrumented to abort if λ is violated. The instrumentation guarantees that λ holds for the program P' , so ϕ holds for P' as well (as long as ϕ is not related to termination properties of P'). Furthermore, if λ holds for P , then P and P' are semantically equivalent. We have in mind security policies, for example the policy ϕ_{flow} that no malicious explicit information flow occurs. The procedure is summarized in Algorithm 1, and an overview of the system (for callgraph specifications discussed in Section 5) is shown in Figure 2.

In our example in Figure 1, we instrument P_{onCreate} to ensure that the cut $\lambda_3 = (3.\text{sendHTTP}(\text{data}) \notin S^*)$ holds. Then `INSTRUMENT(P_{onCreate} , λ_3)` adds instrumentation that terminates P_{onCreate} if `3.sendHTTP(data)` is reached.

3.3 Improving Precision Using Multiple Cuts

We can improve precision by computing multiple sufficient cuts, which must *all* fail before the instrumentation terminates P' . In other words, we want $\lambda_1, \dots, \lambda_n$ such that

$$\chi \wedge (\lambda_1 \vee \dots \vee \lambda_n) \models \phi \text{ and } \forall i, \text{SAT}(\chi \wedge \lambda_i).$$

However, we need to avoid choosing $\lambda_1 = \dots = \lambda_n$ (since then failures are correlated). To do so, we assume that Λ comes with a meet operator \sqcap , where $\lambda \sqcap \mu$ should mean “intersection of specifications λ and μ ”. We require that the predicates be disjoint, i.e. $\lambda_i \sqcap \lambda_j = \perp$ for all $1 \leq i < j \leq n$.

Algorithm 1 Algorithm for interactively verifying P . Here, the function `CUT` solves the abductive inference problem (described in Section 5), and the function `EXTRACTFACTS` constructs the must-facts χ (described in Section 4).

```

procedure INTERACTIVECUT( $P, \phi$ )
   $T \leftarrow \emptyset$ 
  while true do
     $[\chi, \Lambda] \leftarrow \text{EXTRACTFACTS}(P, T)$ 
     $\lambda \leftarrow \text{CUT}(\phi, \chi, \Lambda)$ 
    if  $\lambda = \emptyset$  then
      return  $\emptyset$ 
    end if
     $T_{\text{new}} \leftarrow \mathcal{O}(P, \lambda)$ 
    if  $T_{\text{new}} = \emptyset$  then
      return  $\lambda$ 
    end if
     $T \leftarrow T \cup T_{\text{new}}$ 
  end while
end procedure

procedure INTERACTIVEVERIFY( $P, \phi$ )
   $\lambda \leftarrow \text{INTERACTIVECUT}(P, \phi)$ 
  return INSTRUMENT( $P, \lambda$ )
end procedure

```

We incrementally construct these predicates. Our first predicate is $\lambda_1 \leftarrow \text{CUT}(\phi, \chi, \Lambda)$. When computing λ_2 , we need to ensure that $\lambda_1 \sqcap \lambda_2 = \perp$ —i.e., we need to exclude every predicate in the *downward closure* $(\downarrow \{\lambda_1\}) = \{\mu \in \Lambda \mid \mu \leq \lambda_1\}$ of λ_1 from consideration. To exclude these predicates, we add them to χ :

$$\chi_1 \leftarrow \chi \wedge \bigwedge_{\lambda \in (\downarrow \{\lambda_1\}) - \{\perp\}} (\neg \lambda).$$

Now consider $\lambda_2 \leftarrow \text{CUT}(\phi, \chi, \Lambda)$. Let $\nu = \lambda_1 \sqcap \lambda_2$. Note that $\nu \in (\downarrow \{\lambda_1\})$, since $\nu \leq \lambda_1$. However, `CUT` returns λ_2 such that $\text{SAT}(\chi_1 \wedge \lambda_2)$, and $\chi_1 = (\neg \nu) \wedge (\dots)$ unless $\nu \notin (\downarrow \{\lambda_1\}) - \{\perp\}$, so it must be the case that $\nu = \perp$.

In general, after computing the first $i - 1$ predicates $\{\lambda_1, \dots, \lambda_{i-1}\}$, we compute

$$\chi_i \leftarrow \chi \wedge \bigwedge_{k=1}^{i-1} \bigwedge_{\mu \in (\downarrow \{\lambda_k\}) - \{\perp\}} (\neg \mu),$$

and choose $\lambda_i \leftarrow \text{CUT}(\phi, \chi_i, \Lambda)$. Algorithm 2 uses this procedure to compute $\alpha = \lambda_1 \vee \dots \vee \lambda_n$. Note that at some point, the problem of computing `CUT` becomes infeasible, after which no new sufficient cuts can be computed.

To instrument the program to enforce $\alpha = \lambda_1 \vee \dots \vee \lambda_n$, we keep a global array of Boolean variables $[b_1, \dots, b_n]$, all initialized to false. Whenever a predicate λ_i is violated, we update $b_i \leftarrow \text{true}$. If $b_1 \wedge \dots \wedge b_n$ ever becomes true, then all of the predicates λ_i have been violated and we terminate P' .

Algorithm 2 Algorithm for computing multiple cuts.

```

procedure MULTIPLECUT( $\phi, \chi, \Lambda, n$ )
   $\alpha \leftarrow \text{false}$ 
  for all  $1 \leq k \leq n$  do
     $\lambda_i \leftarrow \text{CUT}(\phi, \chi, \Lambda)$ 
    if  $\lambda_i = \emptyset$  then
       $\alpha \leftarrow \alpha \vee \lambda_i$ 
       $\chi \leftarrow \chi \wedge \bigwedge_{\mu \in (\downarrow \{\lambda_i\}) - \{\perp\}} (\neg \mu)$ 
    end if
  end for
  return  $\alpha$ 
end procedure

```

4. Background on CFL Reachability

We compute cuts for a static (explicit) information flow analysis formulated as a context-free language (CFL) reachability problem. Our analysis does not handle implicit information flows (i.e., where information is leaked due to control flow decisions), and furthermore is flow-, context-, and path-insensitive.

The first step of performing the analysis is to extract a graph $G = (V, E)$ from the given program P , along with a context-free grammar (CFG) $C = (\Sigma, U, P, T)$, where Σ is the set of terminals, U is the set of non-terminals, P is the set of productions, and T is the start symbol. We assume that C is *normalized*, i.e. every production has the form $A \rightarrow B$ or $A \rightarrow BD$ [31]. The edges $e \in G$ are labeled with terminals $\sigma \in \Sigma$, i.e., $e = v \xrightarrow{\sigma} v'$ (for some $v, v' \in V$). The *transitive closure* G^C of G under C is the minimal sized solution to the following constraint system:

- $\frac{e \in G}{e \in G^C}$
- $\frac{v \xrightarrow{B} v' \in G^C \quad A \rightarrow B \in C}{v \xrightarrow{A} v' \in G^C}$
- $\frac{v \xrightarrow{B} v'' \xrightarrow{D} v' \in G^C \quad A \rightarrow BD \in C}{v \xrightarrow{A} v' \in G^C}$

The policy we are trying to enforce has the form $\phi = e^* \notin G^C$, where $e^* = v_{\text{source}} \xrightarrow{T} v_{\text{sink}}$ for the *source* vertex v_{source} and the *sink* vertex v_{sink} . Our exposition focuses on the case where there is a single source and a single sink, but our techniques generalize to policies ϕ excluding edges between multiple sources and sinks.

4.1 Explicit Information Flow Analysis

Given a program P , the static analysis constructs the labeled graph $G = (V, E)$. Here, $V = \mathcal{M} \cup \mathcal{O} \cup \mathcal{V}$, where \mathcal{M} is the set of methods, \mathcal{O} is the set of abstract objects (which correspond to object allocation sites), and \mathcal{V} is the set of reference variables. Additionally, a distinguished *source* vertex $v_{\text{source}} \in \mathcal{V}$ and a distinguished *sink* vertex $v_{\text{sink}} \in \mathcal{V}$ are annotated by the user—for example, in Figure 1, the return value of `getLocation` is annotated as a source, and the argument `text` of `sendHTTP` is a sink.

1. $v = \text{new } X() \Rightarrow o \xrightarrow{\text{New}} v$
2. $u = v \Rightarrow v \xrightarrow{\text{Assign}} u$
3. $u.f = v \Rightarrow v \xrightarrow{\text{Put}[f]} u$
4. $u = v.f \Rightarrow v \xrightarrow{\text{Get}[f]} u$
5. $@\text{Source}(v) \Rightarrow v \xrightarrow{\text{SrcRef}} v$
6. $@\text{Sink}(v) \Rightarrow v \xrightarrow{\text{RefSink}} v$
7. $v \xrightarrow{\sigma} v' \Rightarrow v' \xrightarrow{\bar{\sigma}} v$ (where $\bar{\bar{\sigma}} = \sigma$)

Figure 3. Program fact extraction rules.

8. $\text{FlowsTo} \rightarrow \text{New}$
9. $\text{FlowsTo} \rightarrow \text{FlowsTo Assign}$
10. $\text{FlowsTo} \rightarrow \text{FlowsTo Put}[f] \text{FlowsTo FlowsTo Get}[f]$
11. $\text{SrcSink} \rightarrow \text{SrcRef FlowsTo FlowsTo RefSink}$
12. $A \rightarrow A_1 \dots A_k \Rightarrow \bar{A} \rightarrow \bar{A}_k \dots \bar{A}_1$ (where $\bar{\bar{A}} = A$)

Figure 4. Productions for C_{flow} .

The edges $e \in E$ encode relations between the variables and abstract objects, where the relation is explained by the label $\sigma \in \Sigma$. Rules for generating G are shown in Figure 3. Rules 1-4 handle statements that are used in the points-to analysis [44]. Rule 1 says that the contents of abstract object o flow to the reference v (more formally, v may *point to* o). Rule 2 encodes the flow when a reference variable v is assigned to another reference variable u . Rules 3 and 4 record the flows induced by field reads (or *gets*) and writes (or *puts*)—note that there is a distinct put/get terminal for each field $f \in \mathcal{F}$ (where \mathcal{F} is the set of fields in the program). Additionally, we handle interprocedural information flow by including (i) assignment edges from the arguments in the caller to the formal parameters of the callee, and (ii) assignments from the formal return values of the callee to the left-hand side of the invocation statement in the caller.

Rules 5-6 formalize the source and sink annotations—Rule 5 adds a self-loop on the source vertex, and Rule 6 adds a self-loop on the sink vertex. Finally, Rule 7 is a technical device that allows us to express paths with “backwards” edges. It introduces a label $\bar{\sigma}$ to represent the reversal of an edge labeled σ . Figure 5 shows the part of the graph extracted from the code in Figure 1 (using the rules in Figure 3) that is relevant to finding the source-sink flow from `getLocation.return` to `sendHTTP.param`.

Information flows through the graph correspond to source-sink paths in the CFG C_{flow} defined as follows:

$$\begin{aligned} \Sigma_{\text{flow}} &= \{\text{New}, \text{Assign}, \text{SrcRef}, \text{RefSink}\} \\ &\quad \cup \{\text{Put}[f], \text{Get}[f] \mid f \in \mathcal{F}\} \\ U_{\text{flow}} &= \{\text{New}, \text{Assign}, \text{FlowsTo}, \text{SrcSink}\} \end{aligned}$$

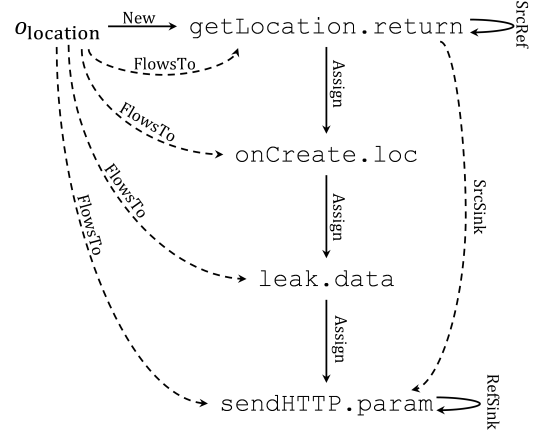


Figure 5. A part of the graph G for the code in Figure 1. Solid edges are edges extracted using the rules in Figure 1. Dashed edges are edges added by the rules in Figure 4. Backwards edges are omitted for clarity.

where \mathcal{F} is the set of fields used in the program. We also include symbols $\bar{\sigma}$ and \bar{A} in Σ_{flow} and U_{flow} , respectively. The start symbol of the grammar is $T_{\text{flow}} = \text{SrcSink}$. Productions are shown in Figure 4. Rules 8-10 build the points-to relation $o \xrightarrow{\text{FlowsTo}} v$, which means that variable $v \in \mathcal{V}$ may point to abstract object $o \in \mathcal{O}$. Rule 11 produces a source-sink edge if the value in source vertex v_{source} is aliased with the value in sink vertex v_{sink} . Finally, Rule 12 introduces a reversed edge $v' \xrightarrow{\bar{A}} v$ for every non-terminal edge $v \xrightarrow{A} v'$. In this way, Rule 12 plays the same role for non-terminal edges that Rule 10 plays for terminal edges. The rules in Figure 4 easily generalized to handle primitive flows. It is also possible to capture implicit flows through framework methods, see [11].

In Figure 5, the dashed edges are edges added when computing the transitive closure. For example, because `return` $\xrightarrow{\text{New}}$ `o_location`, Rules 8 and 12 add the edge `return` $\xrightarrow{\text{FlowsTo}}$ `o_location`. Also, because we have path

$$o_{\text{location}} \xrightarrow{\text{New}} \text{return} \xrightarrow{\text{Assign}} \text{loc} \xrightarrow{\text{Assign}} \text{data} \xrightarrow{\text{Assign}} \text{text},$$

Rules 8 and 9 add edge `o_location` $\xrightarrow{\text{FlowsTo}}$ `text`. Now we have path

$$\text{return} \xrightarrow{\text{SrcRef}} \text{return} \xrightarrow{\text{FlowsTo}} o_{\text{location}} \xrightarrow{\text{FlowsTo}} \text{text} \xrightarrow{\text{RefSink}} \text{text},$$

from which Rule 11 adds `return` $\xrightarrow{\text{SrcSink}}$ `text`.

5. Cuts for CFL Reachability

In this section, we describe an algorithm for performing conditional verification in the case of context-free language reachability. Let C be a context-free grammar, and let $G = (V, E)$ be a labeled graph constructed from a program P . We consider policies ϕ of the form $\phi = (e^* \notin G^C)$, where

$e^* = v_{\text{source}} \xrightarrow{T} v_{\text{sink}}$. This question can be answered in polynomial time [31]. However, the graph constructed by the static analysis in general is an approximation of the true graph $G^* = (V^*, E^*)$, i.e. $G^* \subseteq G$, which potentially introduces false positive source-sink paths. The must-facts include both edges certain to be in G^* (i.e., predicates ($e \in G^*$)), and edges certain *not* to be in G^* (i.e., predicates ($e \notin G^*$)). Our may-facts describe all remaining edges, i.e., edges that may or may not be in G . Let $E_p \subseteq E$ be this set of *may-edges*; then $e \in E_p$ corresponds to may-fact ($e \in G^*$). Our goal is to infer specifications of the form

$$\lambda = \bigwedge_{e \in E_\lambda} (e \notin G^*)$$

where $E_\lambda \subseteq E_p$. In other words, λ specifies that the edges in E_λ are *not* in G^* . The partial ordering on the lattice Λ of specifications is $\lambda \leq \mu$ if $E_\lambda \subseteq E_\mu$, i.e. λ makes fewer assumptions about which edges are not in G^* .

In this setting, the abductive inference problem is to find a minimal subset $E_\lambda \subseteq E_p$ such that $\lambda \models \phi$ holds for P . All else being equal, we prefer to find the smallest cuts possible, so we add the stronger constraint that $|E_\lambda|$ is minimized.

DEFINITION 5.1. Let $G_\lambda = (V, E - E_\lambda)$, i.e. the subgraph of G with edges $e \in E_\lambda$ removed. A predicate $\lambda \in \Lambda$ is a *sufficient cut* if and only if $e^* \notin G_\lambda^C$. The *CFL reachability minimum cut problem* is to find a sufficient cut λ that minimizes $|E_\lambda|$ —i.e., there does not exist any sufficient cut μ such that $|E_\mu| < |E_\lambda|$.

5.1 Algorithms for CFL Reachability Cuts

We describe a reduction of the minimum cut problem to an integer linear program (ILP). The objective of the ILP is to minimize $|E_\lambda|$ over the set of predicates $\{\lambda \in \Lambda \mid e^* \notin G_\lambda^C\}$. We need to translate the constraints on λ into linear inequalities. To do so, we first recast the problem by introducing the Boolean variables $\delta_e = (e \notin G_\lambda^C) \in \{0, 1\}$ for every edge $e \in G^C$ —i.e., $\delta_e = 1$ if removing E_λ from E causes e to be removed from G . We can recover E_λ given the values δ_e , i.e. $E_\lambda = \{e \in E_p \mid \delta_e = 1\}$. In this formulation, the objective is to minimize $|E_\lambda| = \sum_{e \in E_p} \delta_e$.

Recall that $e = v \xrightarrow{A} v' \in G_\lambda^C$ if there exists $e' = v \xrightarrow{B} v''$ and $e'' = v'' \xrightarrow{D} v'$ in G_λ^C such that $A \rightarrow BD \in C$ (we describe the case of binary productions—the case of unary productions is similar); we denote such a triple as $e \rightarrow e'e''$. Then $\delta_e \Rightarrow (\delta_{e'} \vee \delta_{e''})$ must hold—i.e., e is removed from G_λ^C only if either e' or e'' is removed from G_λ^C . Next, for the source-sink edge e^* , we add constraint $\delta_{e^*} = 1$, which enforces ($e^* \notin G_\lambda^C$). Finally, we require that $\delta_e = 0$ for edges $e \in E - E_p$, since these edges cannot be removed from the graph. These constraints translate into linear inequalities:

1. Productions: $\delta_e \leq \delta_{e_1} + \dots + \delta_{e_k}$ for every production $e \rightarrow e_1 \dots e_k$ ($k \in \{1, 2\}$).
2. Remove the source-sink edge: $\delta_{e^*} = 1$.

3. Retain must-edges: $\delta_e = 0$ for every $e \in E - E_p$

The first set of constraints follows because $\delta_e = 1$ only if $\delta_{e_i} = 1$ for some $1 \leq i \leq k$.

The number of constraints generated by this approach is intractable for the typical ILP solver, so we introduce three optimizations to reduce the number of constraints. First, we only need to include variables δ_e for edges $e \in G^C$ —this is because $G_\lambda^C \subseteq G^C$ for every λ , so $\delta_e = 0$ for all $e \notin G^C$.

Second, we construct the constraints in a top-down manner—i.e., we only include productions contained in some derivation of e^* . If an edge $e \in G^C$ is not contained in any derivation of e^* , then the presence of e in G_λ^C does not affect the presence of e^* in G_λ^C , so e can be ignored. This optimization is implemented by first processing all productions $e^* \rightarrow e_1 \dots e_k$ ($k \in \{1, 2\}$); for every input e_i , we recursively add productions for e_i , which recursively adds every production in some derivation of e^* .

Third, any facts added to G^C produced from only the must-edges (i.e., from edges $e \in E - E_p$) will be present in G_λ^C for every $\lambda \in \Lambda$. Note that the graph $G_\top = (V, E - E_p)$ contains no edges $e \in E_p$, so the edges $e \in G_\top^C$ are produced by must-facts alone. This means that we can first compute G_\top^C , and then only include variables δ_e for $e \in (G^C - G_\top^C)$. More precisely, consider a production $e \rightarrow e'e''$:

1. If $e', e'' \in G_\top^C$, then $e \in G_\top^C$, so we do not add any constraints.
2. If $e' \in G_\top^C$ but $e'' \notin G_\top^C$, then we treat this as the *unary* production $e \rightarrow e''$.
3. If $e', e'' \notin G_\top^C$, then we treat this as $e \rightarrow e'e''$ as before.

Algorithm 3 summarizes the procedure. In practice, we include one additional constraint. For $\sigma \in \Sigma$, edges $e_1 = v \xrightarrow{\sigma} v'$ and $e_2 = v' \xrightarrow{\bar{\sigma}} v$ are distinct edges, but they are derived from the same program fact. To account for this, we impose the additional constraint $\delta_{e_1} = \delta_{e_2}$ for such pairs of edges.

THEOREM 5.2. *The set E_p returned by Algorithm 3 solves CFL reachability minimum cut problem.*

For example, consider the graph in Figure 5. Suppose that The graph shown in Figure 6 summarizes the possible derivations of the edge $\text{return} \xrightarrow{\text{SrcSink}} \text{param}$ from the terminal edges; to distinguish this graph from G^C we refer to the edges of this graph as *arrows* and the vertices as *nodes*. There are two types of nodes—nodes corresponding to *productions* $e \rightarrow e_1 \dots e_k$ (shown as black circles), and nodes corresponding to edges in G^C (shown as boxes containing the corresponding edge). Each production $e \rightarrow e_1 \dots e_k$ has one incoming arrow from e , and one outgoing arrow to each of the edges e_1, \dots, e_k . Let

$$E_p = \{v \xrightarrow{\text{Assign}} v' \mid v \text{ formal return value or } v' \text{ formal parameter}\}.$$

In other words, E_p is the set of edges corresponding to method invocations—recall that we treat each method invo-

Algorithm 3 This algorithm solves the CFL reachability minimum cut problem. Here, \mathcal{S} maps variables δ_e to their value in the solution to the ILP.

```

procedure CFLCUT( $C, G, E_p$ )
   $G^C \leftarrow \text{CLOSURE}(C, G)$ ;  $G^C_{\top} \leftarrow \text{CLOSURE}(C, G - E_p)$ 
   $\mathcal{C} \leftarrow \{\delta_{e^*} = 1\}$ 
   $W \leftarrow [e^*]$ ;  $X \leftarrow \{e^*\}$ 
  while  $\neg W.\text{EMPTY}()$  do
     $e \leftarrow W.\text{POP}()$ 
    for all  $e \rightarrow e_1 \dots e_k$  do
       $F \leftarrow \{e_i \mid e_i \notin G^C_{\top}\}$ 
       $\mathcal{C} \leftarrow \mathcal{C} \cup \{\delta_e \leq \sum_{e \in F} \delta_{e_i}\}$ 
       $W \leftarrow W.\text{CONCAT}(\{e \in F \mid e \notin X\})$ 
       $X \leftarrow X \cup \{e \in F \mid e \notin X\}$ 
    end for
  end while
   $\mathcal{S} \leftarrow \text{SOLVEILP}(\min \sum_{e \in E_p} \delta_e, \mathcal{C})$ 
  return  $\{e \in E_p \mid \mathcal{S}(\delta_e) = 1\}$ 
end procedure

```

cation $x = \text{foo}(y)$ as an assignment from argument y to formal parameter foo.param and an assignment from formal return value foo.return to the defined variable x .

Each production generates one constraint $\delta_e \leq \delta_{e_1} + \dots + \delta_{e_k}$ in the ILP, though these constraints are simplified using the three optimizations described above. Figure 7 shows the constraints generated by Algorithm 3. Constraint 1 enforces that the SrcSink edge is in the cut. Constraint 2 enforces the production

$$(\text{return} \xrightarrow{\text{SrcSink}} \text{param}) \rightarrow (\text{return} \xrightarrow{\text{SrcRef}} \text{return} \xrightarrow{\text{FlowsTo}} o_{\text{location}} \xrightarrow{\text{FlowsTo}} \text{param} \xrightarrow{\text{RefSink}} \text{param})$$

but the first, second, and fourth edges are in G^C_{\top} , so they are not included in the constraint. The path $o_{\text{location}} \xrightarrow{\text{FlowsTo}} \text{param}$, which is produced from the three weight 1 edges $\text{return} \xrightarrow{\text{Assign}} \text{loc}$, $\text{loc} \xrightarrow{\text{Assign}} \text{data}$, and $\text{data} \xrightarrow{\text{Assign}} \text{param}$, is captured by Constraints 3-5.

5.2 Multiple CFL Cuts

To compute multiple cuts, we need an efficient way to compute $(\downarrow \{\lambda_i\} - \{\perp\})$. In our case,

$$\bigwedge_{\mu \in (\downarrow \{\lambda_i\} - \{\perp\})} (\neg \mu) = \bigwedge_{e \in E_{\lambda_i}} (\neg \mu_{\{e\}}) = \bigwedge_{e \in E_{\lambda_i}} (e \in G^*).$$

To see the first equality, note that the conjunction on the right-hand side is over a subset of the conjunction on the left-hand side, so the left-hand side implies the right-hand side. Conversely, every $\mu \in (\downarrow \{\lambda_i\} - \{\perp\})$ can be expressed as a (nonempty) conjunction $\mu_{\{e_1\}} \wedge \dots \wedge \mu_{\{e_m\}}$, where $e_1, \dots, e_m \in E_{\lambda_i}$. Therefore $\neg \mu = (\neg \mu_{\{e_1\}}) \vee \dots \vee (\neg \mu_{\{e_m\}})$, which is implied by the right-hand side.

The resulting update rule is $\chi_i \leftarrow \bigwedge_{e \in E_{\lambda_i}} (e \in G^*)$ —i.e., the next call to CUT assumes that every edge $e \in E_{\lambda}$ is in

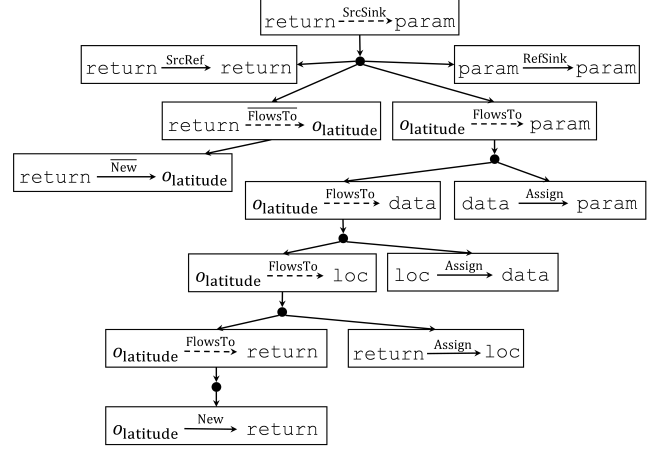


Figure 6. The (partial) derivation tree corresponding to $\text{return} \xrightarrow{\text{SrcSink}} \text{param}$ in Figure 1.

$\max \left\{ \delta(\text{return} \xrightarrow{\text{Assign}} \text{loc}) + \delta(\text{data} \xrightarrow{\text{Assign}} \text{param}) \right\}$ subject to

1. $\delta(\text{return} \xrightarrow{\text{SrcSink}} \text{param}) = 1$
2. $\delta(\text{return} \xrightarrow{\text{SrcSink}} \text{param}) \leq \delta(o_{\text{latitude}} \xrightarrow{\text{FlowsTo}} \text{param})$
3. $\delta(o_{\text{latitude}} \xrightarrow{\text{FlowsTo}} \text{param}) \leq \delta(o_{\text{latitude}} \xrightarrow{\text{FlowsTo}} \text{data}) + \delta(\text{data} \xrightarrow{\text{Assign}} \text{param})$
4. $\delta(o_{\text{latitude}} \xrightarrow{\text{FlowsTo}} \text{data}) \leq \delta(o_{\text{latitude}} \xrightarrow{\text{FlowsTo}} \text{loc}) + \delta(\text{loc} \xrightarrow{\text{Assign}} \text{data})$
5. $\delta(o_{\text{latitude}} \xrightarrow{\text{FlowsTo}} \text{loc}) \leq \delta(o_{\text{latitude}} \xrightarrow{\text{FlowsTo}} \text{return}) + \delta(\text{return} \xrightarrow{\text{Assign}} \text{loc})$

Figure 7. ILP corresponding to productions in Figure 6.

G^* . We define the meet operator as $E_{\lambda_i \sqcap \lambda_j} = E_{\lambda_i} \cap E_{\lambda_j}$ —then, $\lambda_i \sqcap \lambda_j = \perp$ is equivalent to $E_{\lambda_i} \cap E_{\lambda_j} = \emptyset$, which is correctly enforced since χ is updated so that every edge that occurs in E_{λ_i} is prevented from occurring in E_{λ_j} ($j > i$).

6. Implementation

We have implemented the conditional verification algorithm for explicit information flow analysis of Android apps within the Chord program analysis framework [32] modified to use Soot as a front end [48]. We use the ILP solver SCIP [2]. While our information flow analysis (described in Section 4) is not context-sensitive, we compute a 2-CFA points-to analysis in BDDDBDD [50] and use it to filter the points-to set we compute.

App	LOC	FP/TP	Choice of E_p	$ \{e^*\} $	$ E_p $	$ \mathcal{V} $	$ \mathcal{C} $	$ \mathcal{C}_{\text{opt}} $	$\frac{ \mathcal{C}_{\text{opt}} }{ \mathcal{C} }$ (%)	Solve Time (s)	$ E_{\lambda_1} $	$ E_{\lambda_2} $
411524	388744	TP	$E_p^{\text{param}} \cup E_p^{\text{return}}$	4	27687	143875	1982294	264497	13.3	64.350	6	7
0C2B78	322237	TP	E_p^{param}	3	23027	491313	5075599	955606	18.8	29.483	8	10
f7d928	258206	TP	E_p^{param}	4	47517	882370	18969029	1682637	8.9	663.237	11	25
tingshu	239654	TP	E_p^{param}	5	27398	654597	7530383	1280159	17.0	101.259	26	36
16677	199681	TP	$E_p^{\text{param}} \cup E_p^{\text{return}}$	4	40294	423168	6895576	809289	11.7	243.154	4	5
phone	197595	TP	E_p^{param}	3	8334	82596	967559	155905	16.1	11.882	11	11
583cc9	195424	TP	E_p^{param}	4	44565	791567	12575330	1526112	12.1	276.256	20	23
da8c48	189798	TP	E_p^{param}	1	4950	6115	165696	8486	5.1	0.224	1	1
4292c1	155003	TP	$E_p^{\text{param}} \cup E_p^{\text{return}}$	3	39612	1257605	10155092	2452606	24.2	92.545	1	1
5127eb	142217	TP	$E_p^{\text{param}} \cup E_p^{\text{return}}$	2	43092	288688	4579063	548378	12.0	426.992	5	5
1c2514	100437	TP	$E_p^{\text{param}} \cup E_p^{\text{return}}$	1	28366	347151	3182405	649367	20.4	181.696	3	4
wifi	98121	TP	$E_p^{\text{param}} \cup E_p^{\text{return}}$	3	30715	579345	6568009	1128865	17.2	281.255	8	16
browser	346447	FP	E_p^{param}	4	51200	668575	13717560	1291659	9.4	32.079	7	19
00714C	248148	FP	E_p^{param}	4	50915	986358	16784262	1921775	11.4	118.056	21	25
highrail	247377	FP	$E_p^{\text{param}} \cup E_p^{\text{return}}$	3	39095	587201	9309892	1130262	12.1	451.537	9	9
flow	130866	FP	$E_p^{\text{param}} \cup E_p^{\text{return}}$	4	31200	408828	4758531	791712	16.6	48.796	11	12
calendar	125391	FP	$E_p^{\text{param}} \cup E_p^{\text{return}}$	4	31173	225562	3915737	430486	11.0	13.233	5	5
19780d	87448	FP	$E_p^{\text{param}} \cup E_p^{\text{return}}$	4	28153	244151	3742474	467183	12.5	894.685	21	22
aab740	86395	FP	$E_p^{\text{param}} \cup E_p^{\text{return}}$	4	27810	240987	3695096	461002	12.5	305.971	21	21
9d1da3	56125	FP	$E_p^{\text{param}} \cup E_p^{\text{return}}$	5	20055	217060	4320641	420351	9.7	16.862	4	4
018ee7	53187	FP	$E_p^{\text{param}} \cup E_p^{\text{return}}$	3	20739	147711	1952313	267751	9.8	9.844	5	5
ca70f4	44210	FP	$E_p^{\text{param}} \cup E_p^{\text{return}}$	3	9575	56722	455609	106353	23.3	10.402	3	4
battery	33222	FP	$E_p^{\text{param}} \cup E_p^{\text{return}}$	3	13666	97735	1075707	185077	17.2	377.113	12	13
7d43c8	26621	FP	$E_p^{\text{param}} \cup E_p^{\text{return}}$	4	8773	33279	368064	60441	16.4	5.599	3	4
Avg.	83195	-	E_p^{param}	2.32	12068	75622	2437440	337579	15.9	42.36	5.85	7.74

Figure 8. Statistics for some of the Android apps used in the experiments: the number of lines of Jimple byte code (LOC), whether the app exhibited a false positive information flow (FP/TP), the number of source-sink edges $|\{e^* = v_{\text{source}} \xrightarrow{T} v_{\text{sink}} \in G^C\}|$, the number of may-edges $|E_p|$, the number of variables $|\mathcal{V}|$ in the ILP, the unoptimized number of constraints $|\mathcal{C}|$ and the optimized number of constraints $|\mathcal{C}_{\text{opt}}|$, the percentage $|\mathcal{C}_{\text{opt}}|$ compared to $|\mathcal{C}|$, the run time of the ILP solver, and the size of the first and second cuts (on the first iteration of our algorithm). Where relevant, we give statistics for the largest ILP solved for the given app. Also, we include the average values over the entire corpus of 78 apps (for the choice E_p^{param}).

7. Experimental Results

We demonstrate the effectiveness of our approach by interactively verifying a corpus of 78 Android apps¹, including battery monitors, games, wallpaper apps, and contact managers. A number of apps in this corpus contain malicious functionalities that leak sensitive information (contact data, GPS location, and the device ID) to the Internet, and we have ground truth on which information is leaked for each app. Our goal is to apply Algorithm 1 to produce apps proven not to leak sensitive information. The security policy is $\phi_{\text{flow}} = v_{\text{source}} \xrightarrow{\text{SrcSink}} v_{\text{sink}} \notin G^{\text{Cflow}}$ (with multiple source vertices v_{source} and sink vertices v_{sink}), where C_{flow} is the explicit information flow analysis described in Section 4.

As described in Section 2.1, we prune the program by removing provably unreachable statements computing information flows. Also, we make worst-case assumptions about program entry points—i.e., we assume that every potential callback is an entry point (recall that a potential callback is any method in the application that overrides a method in the Android library).

We consider cuts consisting of method invocation statements and return statements, since these statements deter-

mine interprocedural reachability. As described in Section 5, this corresponds to choosing $E_p^{\text{param}} \cup E_p^{\text{return}}$, where

$$E_p^{\text{param}} = \{v \xrightarrow{\text{Assign}} v' \mid v' \text{ formal parameter}\}$$

$$E_p^{\text{return}} = \{v \xrightarrow{\text{Assign}} v' \mid v \text{ formal return value}\},$$

In order to scale to some of the largest apps in our corpus, we needed to further restrict our search space of cuts (which possibly generates larger cuts)—for these apps we chose simply E_p^{param} . The cut λ asserts that certain edges in E_p cannot happen. If an edge in E_p^{param} or E_p^{return} is cut, then we add a statement `assert(false)` immediately before the corresponding method invocation statement.

In our first experiment, we simply run our tool on the corpus of apps and give statistics for the cuts we generate. In our second experiment, the goal is to iteratively generate specifications that describe reachable code according to Algorithm 1. We describe both experiments and our results in the following sections.

7.1 Inferring Cuts

We ran our tool on all the apps in our corpus. The results for twelve of the largest apps, along with all apps with false positives, are shown in Figure 8. We computed two cuts

¹ We will make this corpus of apps publicly available with the publication.

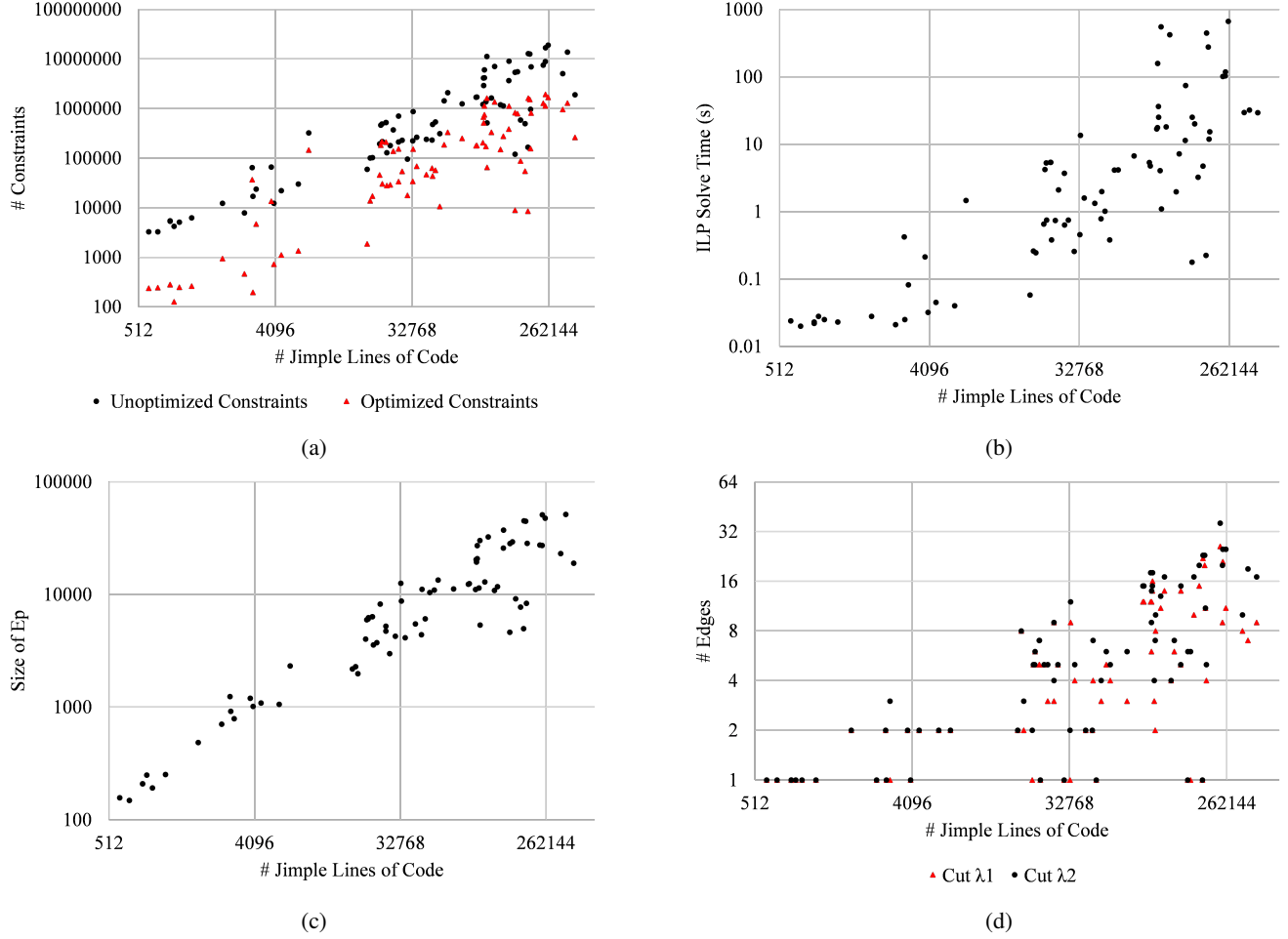


Figure 9. Statistics of the constraint system and resulting cuts for the corpus of 78 Android apps, plotted on a log-log scale: (a) number of unoptimized (black, circle) and optimized (red, triangle) constraints, (b) ILP solve time in seconds, (c) size of the search space E_p , (d) size of the first cut E_{λ_1} (red, triangle) and the second cut E_{λ_2} (black, circle).

for each app, and include the sizes of each of these cuts in Figure 8. In our experience, additional cuts progressively became larger and less useful to examine, though in principle this process can safely be repeated until no new cuts can be produced—having more cuts can only enlarge the set of allowed program behaviors.

We also include some statistics on the sizes of the constraint systems generated by Algorithm 3—these statistics are for the constraint system used to compute the first cut (which is the largest constraint system, though typically the size is similar for other runs). We have shown both the number of unoptimized constraints generated along with the number of constraints after applying the optimizations described in Section 5. Additionally, we include the average values over all 78 apps in the corpus (using E_p^{param} for consistency).

We have plotted some of these statistics in Figure 9 for all the apps in the corpus. For consistency, these figures show results using E_p^{param} . In (a), we compare the size of the un-

optimized constraint system to the size after applying optimizations. As can be seen, the optimized constraint system typically reduces the size by an order of magnitude ($\approx 10\times$). The unoptimized constraint systems typically proved to be intractable for the ILP solver to optimize, but with the optimizations the solver always terminated and finished fairly quickly. We also show the running time for the ILP in (b). As can be seen, our algorithm scales well to apps with hundreds of thousands of lines of Jimple byte code.

In (c), we show the size of E_p^{param} —this gives a sense of the size of the search space of cuts, since there are $2^{|E_p^{\text{param}}|}$ possible cuts. In (d), we show the sizes of the first two cuts produced. Most of the cuts have fewer than 16 edges, though the largest size for the first cut is 26 edges, and the largest size for the second cut is 36 edges. All of these cuts are sufficiently small so that the developer can easily verify whether the cut is valid. This suggests that the interactive verification process places little work on the developer, but we further evaluate this workload in our second experiment.

7.2 Interactive Verification

In our second experiment, we manually carried out the procedure described in Algorithm 1 to produce verified apps P' . Because the app developer is absent, we play the role of the developer.

However, we are disadvantaged compared to the app developer: we only have the app bytecode, have superficial knowledge of the app’s intended functionality, and lack access to the testing tools available to the developer. Furthermore, many of the apps crash when we try to run them due to incompatibilities with the Android emulator. Thus, we provide the reachability information to our tool manually, determining which statements are reachable by reading the Jimple bytecode. For those apps that did not crash in the emulator (about half of the 12 apps) we also ran tests and found that reachability information was consistent with our specifications.

We focus our efforts on producing cuts only for the false positives produced by our explicit information flow analysis (if the flow is a true positive, then no cut exists). For all true positive information flows, the auditor necessarily needs to inspect the app to determine whether it is malicious, so in these cases little can be done to reduce the auditor workload.

The apps with false positive flows are shown in the second half of Figure 8. For each of these apps, we show the source of the false positive flow in Figure 10, and whether we determined that the cause of the false positive is due to unreachable code. These apps typically have other true positive flows—we include only sources that have false positive flows in ϕ_{flow} when performing the verification process (or else ϕ_{flow} would be false for the app).

In Figure 10, we show the results of our interactive verification process. We show two iterations of the process. For each iteration, we show the size of the cuts λ_1 and λ_2 , along with the validity of each cut. The inspection of the cuts proceeded until a valid cut was found, or it was determined that no cut was possible. After just two iterations of Algorithm 1, we succeeded in producing valid cuts for all apps with false positive explicit information flows, except for the app with a false positive not due to unreachable code. This means that only two interactions with the developer were necessary. The cuts remained small after the second iteration, which shows that the entire process is feasible for the developer to carry out. In the case of the final application (browser), because the false positive was not due to unreachable code, no valid cut can be produced by our method. This means that the app would be flagged for manual review.

To demonstrate how each step of Algorithm 1 contributes to verifying each app, Figure 11 plots the number of apps remaining to be verified at each step. As can be seen, the first cut on the first iteration alone clears many of the apps (6 out of 12), and the second cut on the first iteration clears an additional app. The first cut on the second iteration clears three of the remaining apps, and the second cut clears an

additional app, leaving only one app that our process failed to verify.

Whereas the auditor would initially have had to analyze all 12 false positives, our approach reduces the auditor’s workload to a single false positive. In our setting, this may not seem like a huge improvement, because the auditor still needs to analyze the true positive apps. However, our corpus of apps is heavily biased towards apps with malicious behaviors. In practice, the overwhelming majority of apps received by an app store are benign, which means that even a small false positive rate leads to a huge ratio of false positives to true positives that the auditor needs to analyze. We achieve a 92% reduction in the number of false positives that need to be discharged by the auditor, which enables the auditor to better focus effort.

8. Related Work

Our approach to performing iterative verification is related to the following prior techniques.

Abductive inference. Abductive inference has been applied to aid developers in understanding error reports [15], to infer information flow specifications for library functions [56], and for abstraction refinement [53]. We present a general approach for inferring multiple predicates, as well as optimized algorithms for abductive inference in the case where ϕ is described using CFL reachability. Also, [27] presents a related approach to guide sanitizer placement. Their approach is fully automated, but requires runtime taint tracking (though they minimize use of taint tracking).

Specifications from tests. There has been much work on extracting specifications from dynamic executions, with applications to verification [17, 42], proving equivalence of programs [43], finding good abstractions [33], and proving program termination [36]. In another line of work [20], must-facts (extracted from guided dynamic executions) have been used to avoid spending effort trying to discharge true positives. In contrast, we use must-facts extracted from tests as specifications for may-facts—i.e., we *assume* that may-facts not observed in tests are invalid, and use abductive inference to minimize the number of assumptions required.

Dynamic Instrumentation for Safety. Instrumenting programs to ensure safety properties is well-studied, for example to enforce type safety [22, 34] and to ensure control-flow integrity [1]. Our work applies similar principles to ensure the integrity of information flows, which is more challenging because information flows are global properties.

Security applications. Static information flow analysis has been applied previously to the verification of security policies [8, 19, 29, 47, 51]. Our work makes static verification of security policies more practical—rather than employing a large amount of manual labor to discharge false positives, our framework allows the auditor to instrument the program P to enforce the security policy, with the guarantee

App	v_{source}	Cause	Iteration 1				Iteration 2			
			$ E_{\lambda_1} $	$ E_{\lambda_2} $	$\chi \wedge \lambda_1 \stackrel{?}{=} \phi_{\text{flow}}$	$\chi \wedge \lambda_2 \stackrel{?}{=} \phi_{\text{flow}}$	$ E_{\lambda_1} $	$ E_{\lambda_2} $	$\chi \wedge \lambda_1 \stackrel{?}{=} \phi_{\text{flow}}$	$\chi \wedge \lambda_2 \stackrel{?}{=} \phi_{\text{flow}}$
browser	location	unknown	5	6	No	No	5	None	No	No
00714C	contacts	unreachable	2	8	Yes	-	-	-	-	-
highrail	device ID	unreachable	1	2	Yes	-	-	-	-	-
flow	contacts	unreachable	2	3	Yes	-	-	-	-	-
calendar	location	unreachable	3	3	Yes	-	-	-	-	-
19780d	contacts	unreachable	1	1	No	No	2	9	Yes	-
aab740	contacts	unreachable	1	1	No	No	2	9	Yes	-
9d1da3	location	unreachable	4	4	No	No	5	5	No	Yes
018ee7	location	unreachable	4	4	Yes	-	-	-	-	-
battery	location	unreachable	8	8	No	Yes	-	-	-	-
ca7b26	location	unreachable	4	4	Yes	-	-	-	-	-
7d43c8	location	unreachable	3	4	No	No	4	4	Yes	-

Figure 10. Size and validity of cuts generated by Algorithm 3 for apps with false positive flows. “None” means no cut could be generated.

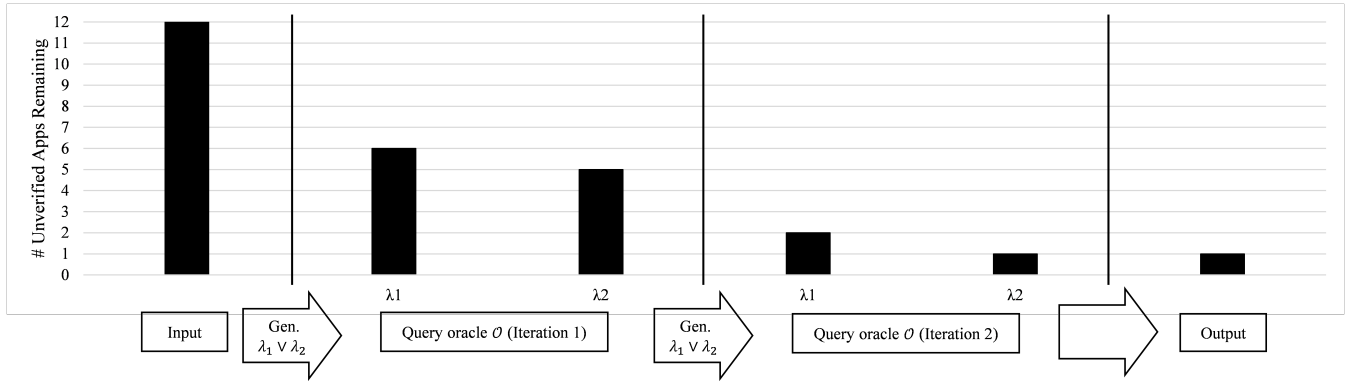


Figure 11. Visualization of how many apps are successfully verified at each step of the process. Algorithm 1 is run on each of the 12 input apps that have a false positive explicit information flow. The x -axis describes the various points in the process, and the y -axis describes the number of apps remaining to be verified at each point.

that the instrumentation is consistent with tests given by the developer.

Dynamic taint tracking has been applied to produce programs that terminate execution upon violation of the security policy [16]. To the best of our knowledge, existing approaches to enforce information flow policies require instrumenting the entire program (or modifying the runtime environment). In contrast, our approach uses very minimal instrumentation, and often places that instrumentation in unreachable code where it will have zero runtime cost. Other approaches for restricting app behaviors have been proposed, for example [40], but the policies enforced are local (e.g., disallowing calls to certain library methods).

Specification inference. There is a large body of work on specification inference; see [11] for a survey. There has also been work specifically on inferring callback specifications [13]. Their approach, which reduces the false positive rate, complements our approach since it would help further reduce the work required of the developer.

9. Conclusions

Given a program P and a policy ϕ , our framework minimally instruments P to ensure that ϕ holds. This instrumentation is guaranteed to be consistent with given test cases, and furthermore the developer can interact with the process to produce suitable cuts. Our approach to handling false positives has the potential to make automated verification of the absence of explicit information flows a more practical approach for security auditors to produce safe and usable programs. We have applied this approach to verify the absence of malicious explicit information flows in a corpus of 78 Android apps. For 11 out of 12 false positives information flows we found, our tool produced valid cuts to enforce ϕ_{flow} .

References

- [1] M. Abadi, M. Budiu, Ú. Erlingsson, J. Ligatti. Control-flow integrity. In *CCS*, 2005.
- [2] T. Achterberg. SCIP: solving constraint integer programs. In *Mathematical Programming Computation*, 2009.
- [3] A. Aiken, S. Bugrara, I. Dillig, T. Dillig, B. Hackett, P. Hawkins. An overview of the Saturn project. In *PASTE*, 43-48,

- 2007.
- [4] R. Alur, P. Černý, P. Madhusudan, W. Nam. Synthesis of interface specifications for Java classes. In *POPL*, 2005.
 - [5] G. Ammons, R. Bodík, J. Larus. Mining specifications. In *POPL*, 2002.
 - [6] Android developers blog, Mar. 2015.
<http://android-developers.blogspot.com/2015/03/creating-better-user-experiences-on.html>
 - [7] Android security blog, Feb. 2012.
<http://googlemobile.blogspot.com/2012/02/android-and-security.html>
 - [8] S. Arzt, S. Rasthofer, C. Fritz, E. Bodden, A. Bartel, J. Klein, Y. L. Traon, D. Octeau, P. McDaniel. FlowDroid: precise context, flow, field, object-sensitive and lifecycle-aware taint analysis for Android apps. In *PLDI*, 2014.
 - [9] N. Ayewah, D. Hovemeyer, J. D. Morgenthaler, J. Penix, W. Pugh. Using static analysis to find bugs. In *IEEE Software*, 2008.
 - [10] T. Ball, S. Rajamani. The SLAM project: debugging system software via static analysis. In *POPL*, 2002.
 - [11] O. Bastani, S. Anand, A. Aiken. Specification inference using context-free language reachability. In *POPL*, 2015.
 - [12] N. Beckman, A. Nori. Probabilistic, modular and scalable inference of typestate specifications. In *PLDI*, 2011.
 - [13] Y. Cao, Y. Fratantonio, A. Bianchi, M. Egele, C. Kruegel, G. Vigna, Y. Chen. EdgeMiner: automatically detecting implicit control flow transitions through the Android framework. In *NDSS*, 2015.
 - [14] A. Chou, B. Chelf, D. Engler, M. Heinrich. Using meta-level compilation to check FLASH protocol code. In *ASPLOS*, 2000.
 - [15] I. Dillig, T. Dillig, A. Aiken. Automated error diagnosis using abductive inference. In *PLDI*, 2012.
 - [16] W. Enck, P. Gilbert, B. Chun, L. Cox, J. Jung, P. McDaniel, A. Sheth. TaintDroid: an information-flow tracking system for realtime privacy monitoring on smartphones. In *OSDI*, 2010.
 - [17] M. D. Ernst, J. H. Perkins, P. J. Guo, S. McCamant, C. Pacheco, M. S. Tschantz, C. Xiao. The Daikon system for dynamic detection of likely invariants. In *Science of Computer Programming*, 2007.
 - [18] Y. Feng, S. Anand, I. Dillig, A. Aiken. Apposcopy: semantics-based detection of Android malware through static analysis. In *FSE*, 2014.
 - [19] A. P. Fuchs, A. Chaudhuri, J. S. Foster. SCanDroid: automated security certification of Android applications. In *IEEE Symposium on Security and Privacy*, 2010.
 - [20] P. Godefroid, A. V. Nori, S. K. Rajamani, S. D. Tetali. Compositional may-must program analysis: unleashing the power of alternation. In *POPL*, 2010.
 - [21] M. Grace, Y. Zhou, Q. Zhang, S. Zou, X. Jiang. RiskRanker: scalable and accurate zero-day Android malware detection. In *MobiSys*, 2012.
 - [22] F. Henglein. Global tagging optimization by type inference. In *ACM Conference on Lisp & Functional Programming*, 1992.
 - [23] J. Kodumal, A. Aiken. Banshee: a scalable constraint-based analysis toolkit. In *SAS*, 2005.
 - [24] J. Kodumal, A. Aiken. Regularly annotated set constraints. In *PLDI*, 2007.
 - [25] J. Kodumal, A. Aiken. The set constraint/CFL reachability connection in practice. In *PLDI*, 2004.
 - [26] T. Kremenek, P. Twohey, G. Back, A. Ng, D. Engler. From uncertainty to belief: inferring the specification within. In *OSDI*, 2006.
 - [27] B. Livshits, S. Chong. Towards fully automated placement of security sanitizers and declassifiers. In *POPL*, 2013.
 - [28] B. Livshits, A. V. Nori, S. K. Rajamani, A. Banerjee. Merlin: specification inference for explicit information flow problems. In *PLDI*, 2009.
 - [29] B. Livshits, M. S. Lam. Finding security vulnerabilities in Java applications with static analysis. In *USENIX Security Symposium*, 2005.
 - [30] B. Livshits, M. S. Lam. Tracking pointers with path and context sensitivity for bug detection in C programs. In *FSE*, 2003.
 - [31] D. Melski, T. Reps. Interconvertibility of a class of set constraints and context-free language reachability. In *Theoretical Computer Science*, 248(1):29-98, 2000.
 - [32] M. Naik, A. Aiken, J. Whaley. Effective static race detection for Java. In *PLDI*, 2006.
 - [33] M. Naik, H. Yang, G. Castelnovo, M. Sagiv. Abstractions from tests. In *POPL*, 2012.
 - [34] G. C. Necula, J. Condit, M. Harren, S. McPeak, W. Weimer. CCured: type-safe retrofitting of legacy software. In *TOPLAS*, 2005.
 - [35] J. W. Nimmer, M. D. Ernst. Automatic generation of program specifications. In *ISSTA*, 2002.
 - [36] A. Nori, R. Sharma. Termination proofs from tests. In *FSE*, 2013.
 - [37] M. K. Ramanathan, A. Grama, S. Jagannathan. Static specification inference using predicate mining. In *PLDI*, 2007.
 - [38] T. Reps. Program analysis via graph reachability. In *ILPS*, 1997.
 - [39] T. Reps, S. Horwitz, M. Sagiv. Precise interprocedural data flow analysis via graph reachability. In *POPL*, 1995.
 - [40] G. Russello, A. B. Jimenez, H. Naderi, W. van der Mark. FireDroid: hardening security in almost-stock Android. In *ACSAC*, 2013.
 - [41] A. Sabelfeld, A. C. Myers. Language-based information-flow security. In *IEEE Journal on Selected Areas in Communications*, 2003.
 - [42] R. Sharma, S. Gupta, B. Hariharan, A. Aiken, A. Nori. Verification as learning geometric concepts. In *SAS*, 2013.
 - [43] R. Sharma, E. Schkufza, B. Churchill, A. Aiken. Data-driven equivalence checking. In *OOPSLA*, 2013.
 - [44] M. Sridharan, D. Gopan, L. Shan, R. Bodik. Demand-driven points-to analysis for Java. In *OOPSLA*, 2005.

- [45] M. Sridharan, S. Artzi, M. Pistoia, S. Guarnieri, O. Tripp, R. Berg. F4F: taint analysis of framework-based web applications. In *OOPSLA*, 2011.
- [46] M. Sridharan, R. Bodik. Refinement-based context-sensitive points-to analysis for Java. In *PLDI*, 2006.
- [47] O. Tripp, M. Pistoia, S. J. Fink, M. Sridharan, O. Weisman. TAJ: effective taint analysis of web applications. In *PLDI*, 2009.
- [48] R. Vallée-Rai, P. Co, E. Gagnon, L. Hendren, P. Lam, V. Sundaresan. Soot: a Java bytecode optimization framework. In *CASCON*, 1999.
- [49] T. Vidas, N. Cristin. Evading Android runtime analysis via sandbox detection. In *ASIA CCS*, 2014.
- [50] J. Whaley, M. S. Lam. Cloning-based context-sensitive pointer alias analysis using binary decision diagrams. In *OOPSLA*, 2004.
- [51] Y. Xie, A. Aiken. Static detection of security vulnerabilities in scripting languages. In *USENIX Security Symposium*, 2006.
- [52] J. Yang, D. Evans, D. Bhardwaj, T. Bhat, M. Das. Perracotta: mining temporal API rules from imperfect traces. In *ICSE*, 2006.
- [53] X. Zhang, R. Mangal, R. Grigore, M. Naik, H. Yang. On abstraction refinement for program analyses in Datalog.
- [54] Y. Zhou, X. Jiang. Dissecting Android malware: characterization and evolution. In *IEEE Symposium on Security and Privacy*, 2012.
- [55] Y. Zhou, Z. Wang, W. Zhou, X. Jiang. Hey, you, get off of my market: detecting malicious apps in official and alternative Android markets. In *NDSS*, 2012.
- [56] H. Zhu, T. Dillig, I. Dillig. Automated inference of library specifications for source-sink property verification. In *APLAS*, 2013.