

CFL-Reachability and Context-sensitive Integrity Types

Ana Milanova Wei Huang Yao Dong

Rensselaer Polytechnic Institute
Troy, NY, USA
{milanova, huangw5, dongy6}@cs.rpi.edu

Abstract

Integrity types can help detect information flow vulnerabilities in web applications and Android apps. We study DFlow, a context-sensitive integrity type system and we give an interpretation of DFlow in terms of CFL-reachability. We propose DFlowCFL, a new, more precise integrity type system, and DFlowCFL-Infer, the corresponding type inference analysis, which is equivalent to CFL-reachability. DFlowCFL-Infer is an effective taint analysis for Android. It scales well and detects numerous privacy leaks in popular Android apps.

Categories and Subject Descriptors D.3.3 [Programming Languages]: Language Constructs and Features

General Terms Experimentation, Languages, Theory

Keywords CFL-Reachability, context-sensitivity, type system, inference

1. Introduction

Pluggable types enhance traditional types. They help find bugs or verify the absence of bugs. With the addition of JSR 308 [7] to Java 8, pluggable types will be increasingly important.

One important class of pluggable types, which we term *information flow types*, helps support information flow control. The type qualifiers and subtyping hierarchy are:

neg <: poly <: pos

Here *neg* is the “negative” qualifier and *pos* is the “positive” qualifier. The goal of the type system is to ensure that there is no flow from a “positive” *source* x to a “negative” *sink* y . *poly* is a polymorphic qualifier, which is interpreted as *pos* in some contexts, and as *neg* in other contexts.

There are many examples of information flow types. Reference immutability types [14, 33] prevent flow from a readonly variable x (readonly is the “positive” qualifier) to a mutable y (mutable is the “negative” qualifier, e.g., y is mutable in $y.f = z$). This ensures that a readonly reference x is never used to mutate the referenced object. Energy types [29] prevent flow from an approximate variable x , which can be represented energy-efficiently, to a precise variable y .

One class of information flow types, which we term *integrity types*, tackles traditional integrity and confidentiality violations [2, 8, 15, 16, 31]. Integrity types prevent flow from untrusted tainted sources (tainted is the “positive” qualifier) to sensitive safe sinks (safe is the “negative” qualifier). They can help uncover violations such as SQL-injection and cross-site scripting (XSS), or verify the absence of such violations. Furthermore, integrity types can help uncover leaks of private data (e.g., phone and location data) to untrusted parties (e.g., Internet) in Android apps, or they can verify the absence of such leaks.¹ For the rest of the paper, we focus on integrity types, although our techniques apply to information flow types in general.

Intuitively, the integrity type system propagates safe sinks backwards through the program. If safe propagates to a variable marked as tainted, then there is a *type error* signaling potential flow from a tainted source to a safe sink; otherwise, the type system guarantees that there is no flow from a source to a sink. Just as with traditional types, integrity types require annotations, which may hinder practical adoption. Therefore, type inference is important.

We advocate *type-based integrity analysis*, which consists of an integrity type system and the corresponding type inference. This analysis is *modular* and *compositional*. It can analyze any given set of classes P . Unknown callees in P are handled using appropriate defaults. Callers of P can be analyzed separately and composed with P without reanalysis of P . Our type-based approach handles reflective object creation seamlessly. This is possible because the type-based analysis does not require abstraction of heap objects; instead, it models flow from one variable to another through subtyping. The analysis requires annotations only on *sources* and *sinks*. Once the sources and sinks are built into annotated libraries, code is analyzed *without* any input from the user. There is evidence that our type-based approach outperforms traditional points-to and dataflow-based integrity analysis [15, 16]. Given these advantages and the advance of JSR 308, it is important to study type-based integrity analysis and its connection to established program analysis techniques.

In this paper, we build upon our recent work on the context-sensitive integrity type system DFlow [15, 16]. We study the connection of DFlow to Context-free Language (CFL)-reachability, which is a well-established program analysis technique [26, 27]. We proceed to define DFlowCFL, a novel more precise integrity type system, and DFlowCFL-Infer, an inference analysis equivalent to CFL-reachability.

The paper makes these concrete contributions:

- CFL-Solver, a novel, CFL-reachability-based integrity analysis.

¹“Integrity types” is clearly a misnomer because in Android the concern is confidentiality. “Confidentiality types” and qualifiers *secret* and *public* would have been more appropriate terms. We keep the term “integrity types” and qualifiers *tainted* and *safe* in deference to previous work, which consistently has used the term “taint analysis” [3, 6].

- DFlowCFL, a new context-sensitive integrity type system.
- Interpretation of DFlow and DFlowCFL in terms of CFL-reachability.
- DFlowCFL-Infer, an inference analysis equivalent to CFL-reachability. DFlowCFL-Infer is an effective taint analysis for Android.

The rest of the paper is organized as follows. Sect. 2 motivates the work. Sect. 3 presents the CFL-reachability-based integrity analysis. Sect. 4 presents the DFlow type system. Sect. 5 presents DFlowCFL and DFlowCFL-Infer, the corresponding inference analysis. Sect. 6 explains the handling of inheritance and virtual calls. Sect. 7 presents a detailed example of the application of DFlowCFL-Infer to the detection of privacy leaks in Android apps. Sect. 8 presents related work and Sect. 9 concludes.

2. Motivation

This work arose from our research on type-based integrity analysis of Java Web applications [15] and Android apps [16].

Static taint analysis for Android especially, is a compelling application of DFlow, DFlowCFL and DFlowCFL-Infer. Taint analysis for Android detects leaks of sensitive data (e.g., phone data, location data) to untrusted places (e.g., the Internet, log files). Static taint analysis for Android has received considerable attention during the last 3 years [3, 8, 10, 12, 17, 19, 20, 36]. One notable effort is DARPA’s Automated Program Analysis for Cybersecurity (APAC) program, which has been active for 2 years. Yet a solution remains elusive. The state-of-the-art analysis is FlowDroid, a highly-precise context-, flow-, field-, object-sensitive and lifecycle-aware analysis for Android [3]. Unfortunately, FlowDroid is heavyweight and memory-intensive. Interestingly, FlowDroid does not uncover malicious leaks in apps from the Google Play Store.

Meanwhile, we worked on a type inference and checking framework, which we instantiated with several type systems and their corresponding inferences [13–15]. DFlow and DFlowCFL-Infer were built as an instance of the framework and incidentally, DFlowCFL-Infer came to be an effective taint analysis for Android. DFlowCFL-Infer scales well, running within a memory footprint of 2GB. It detects many interesting leaks in popular Android apps from the Google Play Store. The detected leaks are interesting because they are consistent (to the best of our understanding) with the widespread belief that “... apps turn smartphones into tracking devices” [1]. The analysis, the handling of Android-specific features, and the empirical evaluation are presented in [16]. (Albeit we call DFlowCFL-Infer DroidInfer in reference [16].) DFlowCFL-Infer’s success prompted our investigation into its connection to CFL-reachability and the discovery of CFL-Solver and DFlowCFL.

The connection to CFL-reachability is important for several reasons. First, it can help address error reporting, a known issue with type-based approaches: given a type error, what source-sink path is to blame for the type error? Our current approach is effective but requires some manual effort and has false positives [16]. With CFL-reachability, we envision error reports with low false-positive rate, where each type error comes with one or more CFL-reachability flow paths from source to sink. Second, as mentioned in the introduction, our type-based analysis obviates the need for pointer analysis as it models flow and aliasing through subtyping. This leads to improved scalability. (We conjecture that the primary reason why DFlowCFL-Infer scales better than FlowDroid, is the fact that DFlowCFL-Infer does not require pointer analysis.) In the same time our type-based analysis is context-sensitive and it achieves good precision [16]. The interpretation of type-based analysis in terms of CFL-reachability, which is one of the most precise flow analysis techniques, can lead

$$\begin{array}{ll}
 C ::= P N \mid P \mid N & F ::= P N \mid P \mid N \\
 P ::= M \mid \rangle_i \mid M P \mid \rangle_i P & P ::= M \mid \rangle_f \mid M P \mid \rangle_f P \\
 N ::= M \mid \langle_i \mid M N \mid \langle_i N & N ::= M \mid \langle_f \mid M N \mid \langle_f N \\
 M ::= d \mid \langle_i M \rangle_i \mid M M & M ::= d \mid \langle_f M \rangle_f \mid M M
 \end{array}$$

(a) Context-free grammar C (b) Context-free grammar F

Figure 1. (a) C handles call-transmitted dependences. i ranges over all callsites. (b) F handles structure-transmitted dependences. f ranges over all fields. $L(C)$ and $L(F)$ denote the languages generated by C and F respectively.

to the development of other scalable and precise type-based flow analyses, beyond taint analysis.

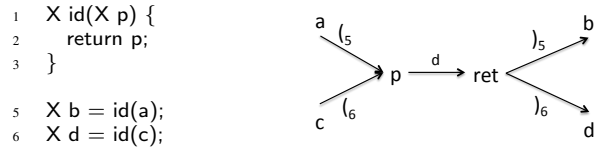
3. CFL-reachability-based Integrity Analysis

This section begins with an overview of CFL-reachability (Sect. 3.1). It proceeds to define the CFL-reachability-based integrity analysis (Sect. 3.2 and Sect. 3.3).

3.1 Overview of CFL-reachability

Context Free Language (CFL)-reachability is an established program analysis technique [26–28]. CFL-reachability casts the problem as a reachability problem over a directed graph G . The nodes in the graph correspond to program variables and the edges represent the flows between these variables.

The following is a classic example of CFL-reachability:



Intraprocedural (direct) flows are represented with d -annotated edges. In the above example, the edge from p to ret represents that parameter p flows to return value ret inside id . For the remainder of the paper, we treat return statements as assignments to a special variable ret (e.g., return p is $ret = p$). Flow from an actual argument to a formal parameter is represented with an \langle_i -annotated edge, where i is the callsite index. The edge from a to p represents the fact that at callsite 5, a flows to p . Finally, flow from the return value to the left-hand-side of the call assignment is represented with \rangle_i -annotated edges, where, again, i is the corresponding callsite.

CFL-reachability casts the problem as follows: if there is a path from x to y in G , such that the annotations on this path form a string in the context-free language of balanced parentheses, then x flows to y . If the annotations do not form a string in this language, x does not flow to y . In the id example above, a flows to b because the string $\langle_5 d \rangle_5$ belongs to the language. However, a does not flow to d because $\langle_5 d \rangle_6$ does not belong to the language. The intuition is clear: CFL-reachability avoids spurious flow across different contexts of invocation, because it matches each call with the corresponding return. The above example illustrates the handling of *call-transmitted dependences* [28] (i.e., dependences due to flow through method calls).

Fig. 1(a) shows grammar C , which handles call-transmitted dependences. M -paths are paths of balanced (matched) parentheses that represent intraprocedural flows. In the id example, d and $\langle_5 d \rangle_5$ are M -paths. N -paths represent paths with outstanding calls, e.g., the path from a to ret has string $\langle_5 d$, which is generated by N . P -paths represent paths due to returns, e.g., the path from ret to b has string \rangle_5 generated by P . We abuse notation slightly by having C denote both the grammar and the starting nonterminal.

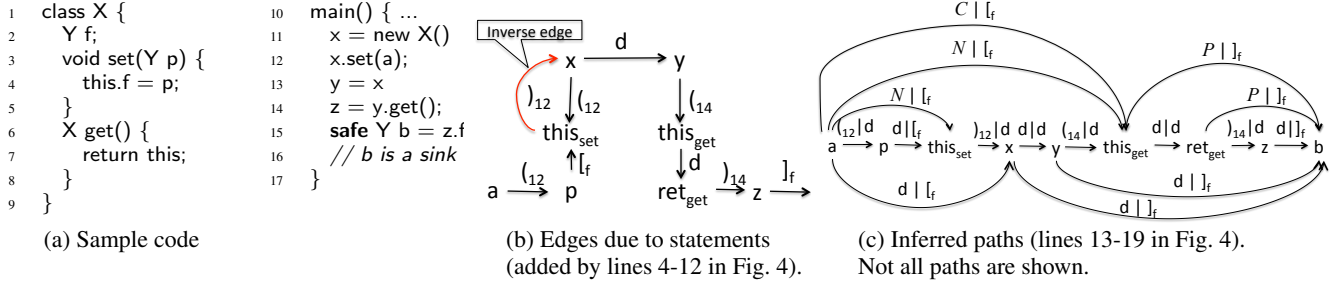
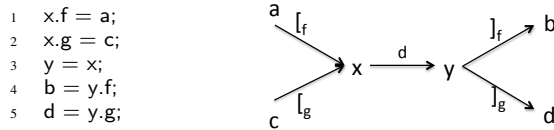


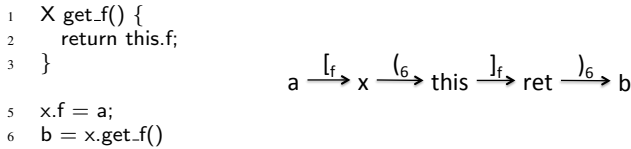
Figure 2. An example with mutable references and interleaving of parentheses and brackets.

Java requires handling of *structure-transmitted dependences* [28] (i.e., dependences due to flow through fields):



Above, $[f]$ -annotated edges model field writes and $]f$ -annotated edges model field reads. The edge from a to x represents the fact that a flows to field f of x (due to $x.f = a$ at line 1). Structure-transmitted dependences can be modeled using an analogous context-free grammar, shown in Fig. 1(b). Here a flows to b because string $[f d]_f$ belongs to the language generated by the grammar in Fig. 1(b). However, a does not flow to d because $[f d]_g$ does not belong to the language.

Had there been only call-transmitted or only structure-transmitted dependences, CFL-reachability could be solved in time $O(n^3)$, where n is the size of the program [26]. Unfortunately, there could be interleaving of the two kinds of flow dependences:



Field f of x is written before the call to `get.f`, it is retrieved in `get.f` and returned into b , creating the path from a to b with interleaved parentheses and brackets. Reps has shown that the precise handling of interleaved parentheses and brackets, i.e., the precise handling of both call-transmitted and structure-transmitted dependences is undecidable [28].

There are additional complications. So far, we create edges as follows: (1) an edge from right-hand-side y to left-hand-side x at direct assignment $x = y$, (2) an edge from actual to formal, (3) an edge from return value to left-hand-side of call assignment, (4) an edge from y to receiver x at field write $x.f = y$, and (5) an edge from receiver w to z at field read $z = w.f$. Consider the code in Fig. 2(a). The call at line 12 to `set` produces two (12) -annotated edges: $x \xrightarrow{(12)} \text{this}_{\text{set}}$ and $a \xrightarrow{(12)} p$. If we were to use only those edges, the flow from a to b would not be discovered! The problem is that `thisset` is *mutated*. The value written into field f of `thisset` is in effect *returned* to the caller of `set`; the “return” must be taken into account.

To overcome this problem we make use of *reference immutability* [14, 33]. Specifically, if the left-hand side of an assignment (explicit or implicit) is *immutable*, then there is one edge, in the expected direction. Otherwise, there is an edge in the expected direction and an additional *inverse edge* in the opposite direction, reversing the annotation as well.

| | |
|--|-----------------------------|
| $cd ::= \text{class } C \text{ extends } D \{ \overline{fd} \overline{md} \}$ | <i>class</i> |
| $fd ::= t' f$ | <i>field</i> |
| $md ::= t m(t \text{ this}, t x) \{ \overline{t} \overline{y} s; \text{return } y \}$ | <i>method</i> |
| $s ::= s; s \mid x = \text{new } t() \mid x = y \mid x = y.f \mid y.f = x \mid x = y.m(z)$ | <i>statement</i> |
| $t ::= q C$ | <i>qualified type</i> |
| $t' ::= q' C$ | <i>qualified field type</i> |
| $q ::= \text{tainted} \mid \text{poly} \mid \text{safe}$ | <i>qualifier</i> |
| $q' ::= \text{tainted} \mid \text{poly}$ | <i>field qualifier</i> |

Figure 3. Syntax. C and D are class names, f is a field name, m is a method name, and $x, y,$ and z are names of local variables, formal parameters, or parameter `this`. As in some of the code examples, `this` is explicit. For simplicity, we assume all names are unique.

At line 12 in Fig. 2, there is one edge due to the implicit assignment of actual a to formal p : $a \xrightarrow{(12)} p$. This edge suffices because p is immutable. However, there are two edges due to the implicit assignment of x to `thisset`: $x \xrightarrow{(12)} \text{this}_{\text{set}}$ and inverse edge $\text{this}_{\text{set}} \xrightarrow{(12)} x$. Note that the annotation is reversed as well, (12) becomes $]12$; this expresses the fact that the value written into f of `thisset` is *returned* to the caller of `set`. The inverse edge is added because there is mutation of `thisset` at `this.f = p`. With the inverse edge, there is a path from a to b with string $(12 [f]]12 d (14 d)]14]f$ as shown in Fig. 2(b). Parentheses and brackets both match.

The intuition behind inverse edges is the following. If there is a write into field f of some y , each path $x \rightarrow \dots \rightarrow y$ must be reversed, because field f can be read out of x , or out of any z to which x flows. Consider

```

1 w = x;
2 y = x;
3 y.f = z;
4 v = w.f;

```

In this example, slightly modified from [32], $y.f$ and $w.f$ are aliases and z flows to v . The inverse edge $y \rightarrow x$ is added because y is modified at $y.f = z$. This edge connects y to x , and then to w , which reveals the flow from z to v .

Inverse edges and inverse paths are similar to the ones in Sridharan and Bodik [32]. The key novelty in our work is the use of reference immutability, which obviates the need for heap abstraction. Sridharan and Bodik [32] require heap abstraction and hence they incur the notorious difficulties of dealing with reflective object creation.

3.2 Algorithm

This section formalizes the intuition into a CFL-reachability-based integrity analysis. As noted, the precise handling of both call-transmitted and structure-transmitted dependences is undecidable,

```

1: procedure CFL-SOLVER
2: Initialize  $G$  with sink node  $n$ 
3: repeat
4:   for each (ASSIGN)  $x = y$  do BIEDGE( $x, y, x, d, d$ ) end for
5:   for each (WRITE)  $y.f = x$  do BIEDGE( $y.f, x, y, d, ]_f$ ) end for
6:   for each (READ)  $x = y.f$  do BIEDGE( $x, y, x, d, ]_f$ ) end for
7:   for each (CALL)  $i : x = y.m(z)$  do
8:     let  $m : \text{this}, p \rightarrow \text{ret}$  be the static target at  $i$ 
9:     BIEDGE( $\text{this at } i, y, \text{this}, (i, d)$ )
10:    BIEDGE( $p$  at  $i, z, p, (i, d)$ )
11:    BIEDGE( $x, \text{ret}, x, (i, d)$ )
12:   end for
13:   for each  $x \xrightarrow{d} y \xrightarrow{d} z \in G$  do EDGE( $x, z, d, f \oplus f'$ ) end for
14:   for each  $x \xrightarrow{(i)} p \xrightarrow{d} p' \xrightarrow{i} z \in G$  do EDGE( $x, z, d, f$ ) end for
15:   for each  $x \xrightarrow{Y} y \in G$  s.t.  $X ::= Y$  do EDGE( $x, y, X, f$ ) end for
16:   for each  $x \xrightarrow{Y} y \xrightarrow{Z} z \in G$  s.t.  $X ::= YZ$  do
17:     EDGE( $x, z, X, f \oplus f'$ )
18:   end for
19:   until  $G$  remains unchanged
20: end procedure
21: procedure BIEDGE( $lhs, s, t, c, f$ )
22:   EDGE( $s, t, c, f$ )
23:   if  $ReIm(lhs) \neq \text{readonly}$  then EDGE( $t, s, \bar{c}, \bar{f}$ ) end if
24: end procedure
25: procedure EDGE( $s, t, c, f$ )
26:   if  $t \xrightarrow{c} n \in G \wedge f \in L(F) \wedge f$  is  $[_f \Rightarrow ]_f \in G$  then
27:     Add  $s \xrightarrow{c} t$  to  $G$ 
28:   end if
29: end procedure

```

Figure 4. CFL-reachability-based integrity analysis. $x \xrightarrow{c} y \in G$ tests if the edge is in G , and $]_f \in G$ tests if there is an edge in G with field annotation $]_f$.

and therefore our analysis approximates. The exact approximation is explained in Sect. 3.3.

We define the analysis over a syntax in “named form” where the results of field accesses, method calls, and instantiations are immediately stored in a variable. The syntax is shown in Fig. 3. Without loss of generality, we assume that methods have parameter this , and exactly one other formal parameter. The type qualifiers present in the figure become relevant in Sect. 4.

The analysis separates the grammar of call-transmitted dependences from the grammar of structure-transmitted dependences. This is achieved by having each edge in the graph present with two orthogonal annotations, one that tracks call-transmitted dependences and another one that tracks structure-transmitted dependences. We write $x \xrightarrow{c} y$. c is the *call annotation*; it is a string in $L(C)$, or a nonterminal in C representing a string in $L(C)$. f is the *field annotation*; it is a string in $L(F)$, or a nonterminal in F .

The analysis, called CFL-Solver, is shown in Fig. 4. It takes as input a set of classes P and a sink variable n . (Without loss of generality we assume that there is a single sink n .) CFL-Solver propagates n through P ; it builds the graph “backwards”, i.e., in opposite direction of the flow edges, adding nodes and edges as it discovers paths to n . CFL-Solver outputs a directed graph G that safely approximates flows to n . If there is flow from reference x to n , then x is in G and n is reachable from x . If x is tainted, CFL-Solver signals potential flow from a source to a sink. For clarity, at this point we assume that every virtual call has a single target. We discuss inheritance and virtual calls in Sect. 6.

CFL-Solver makes use of ReIm, a type system for reference immutability we developed in previous work [14]. ReIm comes with

a quadratic type inference analysis. If ReIm infers that a reference x is readonly, this guarantees that x is never used to mutate the referenced object nor anything it references. If x is readonly, all of the following are forbidden:

- $x.f = z$
- $x.\text{setField}(z)$ where setField sets a field of its receiver
- $y = \text{id}(x)$; $y.f = z$
- $x.f.g = z$
- $y = x.f$; $y.g = z$

Function $ReIm(expr)$ returns the ReIm type of the expression argument. Procedure BIEDGE(lhs, s, t, c, f), shown at lines 21–24, adds edges due to explicit and implicit assignments. It takes lhs — the left-hand-side of the assignment, s and t — the source and target of the edge respectively, and c and f — the call and field annotations of the new edges. CFL-Solver calls BIEDGE at every program statement (lines 4–12). If the left-hand-side of the assignment is readonly according to ReIm, BIEDGE adds only one edge, from s to t . Otherwise, BIEDGE adds an inverse edge from t to s as well. \bar{c} and \bar{f} are as expected: $\bar{[_f]} =]_f$, $\bar{]}_f] = [_f$, $\bar{(i)} =)i$, $\bar{)}i) = (i)$, $\bar{d} = d$. As an example, consider (WRITE) $y.f = x$. If left-hand-side $y.f$ is readonly, then BIEDGE creates edge $x \xrightarrow{d}]_f y$. Otherwise, it creates edges $x \xrightarrow{d}]_f y$ and $y \xrightarrow{d} [_f x$. Fig. 2(b) shows the edges added due to lines 4–12. (There are details we clarify shortly.)

Lines 13–14 add edges that represent M -paths, that is, paths of balanced parentheses. To reduce clutter we use d instead of M . For example, line 13 takes edge $x \xrightarrow{d} y$ and edge $y \xrightarrow{d} z$ and adds a new edge $x \xrightarrow{d} f \oplus f' z$, where $f \oplus f'$ is the concatenation of the two field annotations f and f' . Lines 15–18 add P -paths, N -paths and C -paths (recall Fig. 1(a)). Consider Fig. 2 and edges $y \xrightarrow{(14)} p \xrightarrow{d} p' \xrightarrow{i} z$. Line 15 adds edge $\text{this}_{\text{get}} \xrightarrow{N} \text{ret}_{\text{get}}$ due to production $N ::= M$. Lines 16–18 then add edge $y \xrightarrow{N} \text{ret}_{\text{get}}$ due to production $N ::= (i N$, recording the N -path from y to ret_{get} .² Fig. 2(c) shows P -paths, N -paths and C -paths.

Note that CFL-Solver collapses the strings in $L(C)$ into representative nonterminals (thus eschewing explicit computation of these strings). However, it explicitly computes strings in $L(F)$ (of course, subject to approximation as explained shortly). This means that we have chosen to handle call-transmitted dependences precisely but approximate in the handling of structure-transmitted dependences.

Procedure EDGE (lines 25–29) adds an edge from source s to target t with call annotation c and field annotation f . It adds the edge only if the following conditions are met: (1) sink n is reachable from target t on a path in $L(C)$, (2) the field annotation f is in $L(F)$ and (3) if f is a field write $[_f$, the corresponding field read $]_f$ must be in G . The role of condition (3) becomes clear in the following subsection.

3.3 Approximation

As stated earlier, the precise handling of both call-transmitted and structure-transmitted dependences is undecidable. Therefore, analyses must approximate in the handling of one or both of these dependences. Analyses typically approximate in the handling of call-transmitted dependences [32].

In contrast to previous work, our analysis handles call-transmitted dependences precisely but approximates in the handling of structure-

² Grammar C can be written into Chomsky normal form, which will render lines 13–14 unnecessary. We separate the handling of M -paths (lines 13–14) from the handling of the rest of the paths (15–18) to better illustrate the parallel with type system DFlow.

transmitted dependences. Recall that CFL-Solver was designed to explain DFlow. DFlow as well as other information flow type systems [14, 29] approximate in the handling of structure-transmitted dependences. Our approximation is simple: $f \oplus f'$ always evaluates to d , thus *erasing* field annotations. Just as any approximation, this approximation may introduce imprecision.

Condition (3) in EDGE mitigates the imprecision. If field f never flows to sink n , then $]_f$ is never added to G and the sink does not propagate towards writes of f . Consider

```

1 x.f = a;
2 x.g = c;
3 y = x;
4 b = y.f;
5 safe d = y.g; // d is sink variable n

```

Only field g is added to the graph, f is not added because b does not flow to the sink. Thus, the safe sink is propagated backwards to c as expected, but it is not propagated to a . Imprecision arises *only* when there are two instances of the same class, say x and y , and two distinct fields, say $x.f$ and $y.g$ flow to the sink. The sink is propagated backwards towards $x.g$ and $y.f$ even though these fields may not flow to the sink. In our experience with both Java web applications and Android apps, this is rarely a source of imprecision, because sources and sinks are sparse.

CFL-Solver is in fact a *framework* for CFL-reachability analysis. Analysis designers can choose a different approximation for fields and encode this approximation through grammar F and concatenation operation \oplus . For example, one can reduce F to a regular grammar just as one can reduce C to a regular grammar [32]. Furthermore, one can handle structure-transmitted dependences precisely but approximate in the handling of call-transmitted dependences. We plan to explore these directions in future work.

From now on, when referring to CFL-Solver, we mean the solver implementing the above approximation for fields.

4. Type-based Integrity Analysis

This section describes the type-based integrity analysis. We outline the DFlow type system and interpret it in terms of the CFL-reachability analysis.

4.1 Type Qualifiers

Each variable is typed by a *type qualifier* (Fig. 3). There are two basic qualifiers in DFlow: tainted and safe.

- **tainted**: A variable x is tainted, if there is flow from a source to x . Sources are tainted.
- **safe**: A variable x is safe if there is flow from x to a sink. Sinks are safe.

In order to disallow flow from tainted sources to safe sinks, DFlow enforces the following subtyping hierarchy:

$$\text{safe} <: \text{tainted}^3$$

where $q_1 <: q_2$ denotes q_1 is a subtype of q_2 (q is also a subtype of itself $q <: q$). Therefore, it is allowed to assign a safe variable to a tainted one:

```

safe String s = ...;
tainted String t = s;

```

However, it is not allowed to assign a tainted variable to a safe one:

```

tainted String t = ...;
safe String s = t; // type error!

```

³This is the desired subtyping. However, it is a well-known issue that subtyping is not always safe when *mutable* references are involved [15, 25, 29]. This is precisely the issue that necessitated inverse edges in CFL-Solver.

$$\begin{array}{c}
\text{(TASSIGN)} \\
\frac{\Gamma(x) = q_x \quad \Gamma(y) = q_y \quad q_y <: q_x \\
\text{ReIm}(x) \neq \text{readonly} \Rightarrow q_x <: q_y}{\Gamma \vdash x = y} \\
\\
\text{(TWRITE)} \\
\frac{\Gamma(y) = q_y \quad \text{typeof}(f) = q_f \quad \Gamma(x) = q_x \quad q_x <: q_y \triangleright q_f \\
\text{ReIm}(y.f) \neq \text{readonly} \Rightarrow q_y \triangleright q_f <: q_x}{\Gamma \vdash y.f = x} \\
\\
\text{(TREAD)} \\
\frac{\Gamma(y) = q_y \quad \text{typeof}(f) = q_f \quad \Gamma(x) = q_x \quad q_y \triangleright q_f <: q_x \\
\text{ReIm}(x) \neq \text{readonly} \Rightarrow q_x <: q_y \triangleright q_f}{\Gamma \vdash x = y.f} \\
\\
\text{(TCALL)} \\
\frac{\text{typeof}(m) = q_{\text{this}}, q_p \rightarrow q_{\text{ret}} \\
\Gamma(y) = q_y \quad \Gamma(x) = q_x \quad \Gamma(z) = q_z \quad \Gamma(i) = q^i \\
q_y <: q^i \triangleright q_{\text{this}} \quad \text{ReIm}(\text{this at } i) \neq \text{readonly} \Rightarrow q^i \triangleright q_{\text{this}} <: q_y \\
q_z <: q^i \triangleright q_p \quad \text{ReIm}(p \text{ at } i) \neq \text{readonly} \Rightarrow q^i \triangleright q_p <: q_z \\
q^i \triangleright q_{\text{ret}} <: q_x \quad \text{ReIm}(x) \neq \text{readonly} \Rightarrow q_x <: q^i \triangleright q_{\text{ret}}}{\Gamma \vdash i : x = y.m(z)}
\end{array}$$

Figure 5. DFlow type system. Function *typeof* retrieves the DFlow types of fields and methods, Γ is a type environment that maps variables to DFlow qualifiers.

4.2 Context Sensitivity

DFlow achieves context sensitivity by using a polymorphic type qualifier, poly, and *viewpoint adaptation* [5].

- **poly**: The poly qualifier expresses context sensitivity (i.e., polymorphism). poly is interpreted as tainted in some contexts and as safe in other contexts.

The subtyping hierarchy becomes

$$\text{safe} <: \text{poly} <: \text{tainted}$$

The concrete value of poly is interpreted by the viewpoint adaptation operation. Viewpoint adaptation of a type q' from the viewpoint of another type q , results in the adapted type q'' . This is written as $q \triangleright q' = q''$. Viewpoint adaptation adapts fields, formal parameters, and method return values from the viewpoint of the *context* at the field access or method call. DFlow defines the viewpoint adaptation operation below:

$$\begin{array}{lcl}
- \triangleright \text{tainted} & = & \text{tainted} \\
- \triangleright \text{safe} & = & \text{safe} \\
q \triangleright \text{poly} & = & q
\end{array}$$

The underscore denotes a “don’t care” value. Qualifiers tainted and safe do not depend on the viewpoint (context). Qualifier poly depends on the viewpoint: e.g., if the viewpoint (context) is tainted, then poly is interpreted as tainted. Viewpoint adaptation is applied at field accesses and method calls; we elaborate upon this in the next section.

4.3 Type System

The DFlow typing rules are defined in Fig. 5. The key to the relation of CFL-Solver to DFlow is that (roughly speaking) if CFL-Solver finds a path from x to y , DFlow guarantees $x <: y$. (For the remainder of the paper, we often write x instead of q_x for the type qualifier of x .)

Rule (T_{ASSIGN}) enforces the expected subtyping constraints. If the left-hand-side of the assignment is not readonly according to ReIm , (T_{ASSIGN}) creates an additional subtyping constraint in the opposite direction, just as CFL-Solver creates an inverse edge. The rules for field access, (T_{TREAD}) and (T_{WRITE}) , adapt field f from the viewpoint of receiver y and then enforce the subtyping constraints.

The rule for method call, (T_{CALL}) , adapts formal parameters this and p and return value ret from the viewpoint of *callsite context* q^i , and creates the subtyping constraints that capture flows from actual arguments to formal parameters, and from return value to the left-hand-side of the call assignment. Again, there is analogy with CFL-Solver — constraint $q_y <: q^i \triangleright q_{\text{this}}$ corresponds to edge $y \xrightarrow{i} \text{this}$ and constraint $q^i \triangleright q_{\text{ret}} <: q_x$ corresponds to edge $\text{ret} \xrightarrow{i} x$.

Let us return to the examples from Sect. 3. Consider

```

1  poly X id(poly X p) {
2    return p;
3  }

5  safe X b = id(a);
6  tainted X d = id(c);

```

Statement $\text{return } p$ enforces subtyping constraint $p <: \text{ret}$. The call to id at 5 creates the following constraints (rule (T_{CALL}) in Fig. 5):

$$a <: q^5 \triangleright p \quad q^5 \triangleright \text{ret} <: b$$

It is easy to show that viewpoint adaptation preserves subtyping: that is, if $p <: \text{ret}$ holds, then $q^5 \triangleright p <: q^5 \triangleright \text{ret}$ holds. This constraint, combined with the above two constraints results in $a <: b$. This mirrors line 14 in CFL-Solver. Constraint $a <: q^5 \triangleright p$ corresponds to the (\triangleright_5) -annotated edge from actual a to formal p , $p <: \text{ret}$ corresponds to the (\triangleright) -annotated edge from p to ret , and finally, $q^5 \triangleright \text{ret} <: b$ corresponds to the (\triangleright_5) -annotated edge from ret to the left-hand-side of the call assignment b .

DFlow is context-sensitive. id is polymorphic:

poly X $\text{id}(\text{poly } X \text{ p})$

At callsite 5, the poly type is instantiated to safe and at callsite 6, it is instantiated to tainted . As discussed, DFlow enforces the following constraints at callsite 5:

$$a <: q^5 \triangleright \text{poly} \quad q^5 \triangleright \text{poly} <: \text{safe}$$

$q^5 = \text{safe}$ satisfied the above constraints. DFlow enforces the following constraints at callsite 6:

$$c <: q^6 \triangleright \text{poly} \quad q^6 \triangleright \text{poly} <: \text{tainted}$$

$q^6 = \text{tainted}$ satisfied the above constraints.

Consider the handling of structure-transmitted dependences in DFlow. Recall Fig. 3. It stipulates that fields are tainted or poly. This mirrors CFL-Solver. If the left-hand-side of a field read flows to sink, the field is forced to poly, “erasing” the field info and propagating safe towards writes of the field; otherwise, the field is tainted and right-hand-sides of writes to the field can remain tainted. Consider:

```

1  x.f = a;
2  x.g = c;
3  y = x;
4  b = y.f;
5  safe d = y.g;

```

At field read $d = y.g$, DFlow enforces constraint $y \triangleright g <: \text{safe}$. Since field g is tainted or poly, but not safe, g must be poly. As a result, constraint $y \triangleright g <: \text{safe}$ entails $y <: \text{safe}$ (because $y \triangleright \text{poly}$ equals y). At $x.g = c$, DFlow creates constraint $c <: x \triangleright g$, which again entails $c <: x$. The assignment $y = x$ enforces $x <: y$. The three constraints

$$c <: x \quad x <: y \quad y <: \text{safe}$$

entail $c <: \text{safe}$, as expected. Note that x and y must be safe. However, field f as well as a and b may remain tainted. Again, DFlow mirrors CFL-Solver: the safe sink d is propagated backwards to c , but a remains unaffected.

As a final example, consider the code in Fig. 2. Due to $b = z.f$, DFlow enforces constraint $z \triangleright f <: \text{safe}$, which forces f to poly and entails constraint $z <: \text{safe}$. Since f is poly, statement $\text{this.f} = p$ in set entails $p <: \text{this}$. The call to set at 12 deserves attention. DFlow creates the following constraints:

$$x <: q^{12} \triangleright \text{this} \quad q^{12} \triangleright \text{this} <: x \quad a <: q^{12} \triangleright p$$

(Due to the mutation of this , DFlow enforces $q^{12} \triangleright \text{this} <: x$. This is analogous to the inverse edge in CFL-Solver.) Constraints

$$a <: q^{12} \triangleright p \quad p <: \text{this} \quad q^{12} \triangleright \text{this} <: x$$

entail $a <: x$. Further, DFlow gives rise to these constraints

$$a <: x \quad x <: y \quad y <: \text{safe}$$

mirroring the path from a to b in Fig. 2. The above constraints force a to be safe, as expected.

The above examples give intuition into DFlow and its connection to CFL-Solver. The theorems below formalize the connection. Lemma 4.1 states (roughly) that if CFL-Solver discovers a path from reference variable x to reference variable y , then DFlow enforces appropriately that x is a subtype of y . For example, if CFL-Solver discovers an N -path from x to y , then DFlow enforces constraint $x <: q \triangleright y$ for some value of q .

LEMMA 4.1. *Let $\text{CFL-SOLVER} \vdash e$ denote that CFL-Solver adds edge e to G . Let $\text{DFLOW} \vdash c$ denote that DFlow enforces constraint c . $q, q' \in \{\text{tainted}, \text{poly}, \text{safe}\}$.*

1. $\text{CFL-SOLVER} \vdash x \xrightarrow{d} y \Rightarrow \text{DFLOW} \vdash x <: y$
2. $\text{CFL-SOLVER} \vdash x \xrightarrow{N} y \Rightarrow \text{DFLOW} \vdash x <: q \triangleright y$
3. $\text{CFL-SOLVER} \vdash x \xrightarrow{P} y \Rightarrow \text{DFLOW} \vdash q \triangleright x <: y$
4. $\text{CFL-SOLVER} \vdash x \xrightarrow{PN} y \Rightarrow \text{DFLOW} \vdash q \triangleright x <: q' \triangleright y$

THEOREM 4.2. *Let $\text{CFL-SOLVER} \vdash e$ denote that edge e is in G and let $\text{DFLOW} \vdash c$ denote that constraint c holds.*

1. $\text{CFL-SOLVER} \vdash x \xrightarrow{d} n \Rightarrow \text{DFLOW} \vdash x <: \text{safe}$
2. $\text{CFL-SOLVER} \vdash x \xrightarrow{N} n \Rightarrow \text{DFLOW} \vdash x <: \text{safe}$
3. $\text{CFL-SOLVER} \vdash x \xrightarrow{P} n \Rightarrow \text{DFLOW} \vdash x <: \text{poly}$
4. $\text{CFL-SOLVER} \vdash x \xrightarrow{PN} n \Rightarrow \text{DFLOW} \vdash x <: \text{poly}$

Theorem 4.2 follows directly from Lemma 4.1. For example, if CFL-Solver discovers an N -path from x to n , then by Lemma 4.1 we have $x <: q \triangleright n$. Since n is safe, $x <: \text{safe}$.

5. A Type System Equivalent to CFL-Solver

DFlow is simple and therefore appealing. Unfortunately, it is stricter than CFL-Solver: it propagates sinks further and rejects programs that CFL-Solver deems safe. Theorem 4.2 reflects this. While a path from x to n in G implies that x is safe or poly in DFlow, x being safe or poly does not necessarily guarantee a path from x to n .

This section elaborates on the imprecision in DFlow (Sect. 5.1) and proceeds to define a type system equivalent to CFL-Solver (Sect. 5.2). Finally, it outlines the type inference analysis (Sect. 5.3).

5.1 Imprecision in DFlow

The imprecision arises from rule (T_{CALL}) . Essentially, (T_{CALL}) may force flow from an actual to the left-hand-side of the call assignment, even when there is no flow from the corresponding formal parameter to the return value. Suppose there is a method m , such that there

```

1 X m(A this) { ...
2   this.f = y; ...
3   X x = new X();
4   return x;
5 }
6 ...
7 x1 = a.m7();
8 x1.g = 0;
9 safe Y sink = a.f;
10 ...
11 safe X sink1 = a1.m11();
12 Y x2 = a1.f;

```

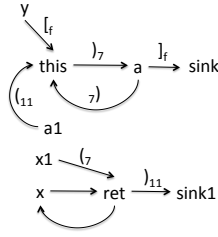


Figure 6. Imprecision in DFlow.

is no flow from its `this` to its `ret`. Suppose that `m`'s `this` flows to `sink` in context q^i , but it does not flow to `sink` in context q^j . (More precisely, it would be `this.f` written in `m` that would flow to `sink`.) Also, let `ret` flow to `sink` in context q^j , but not in context q^i .

Naturally, `m`'s `this` and `ret` should both be poly. Thus `this` $<$: `ret`, which in DFlow leads to

$$y <: q^j \triangleright \text{this} <: q^j \triangleright \text{ret} <: x$$

where `y` and `x` are the receiver and left-hand-side, respectively, at callsite j . Since `x` flows to `sink`, `x` is safe, and the above constraints force `y` to be safe, unnecessarily.

The example in Fig. 6 illustrates this issue and differentiates with CFL-Solver. Reference variables `a` and `a1` refer to two distinct `A` objects. Field `a.f` flows to `sink`, but `x1`, returned out of `a.m()`, does not flow to `sink` (lines 7-9). In contrast, `a1.f` does not flow to `sink`, but `a1.m()` does (lines 11-12). Fig. 6 shows the edges added due to program statements by CFL-Solver. (CFL-Solver proceeds to add path edges accordingly. These edges are not shown.) Note that `a1` and `x1` remain unaffected by the `sink` because there is no C -path from `a1` or `x1` to a `sink`. `x2` is not in the graph.

DFlow enforces the following constraints at line 7:

$$a <: q^7 \triangleright \text{this} \quad q^7 \triangleright \text{this} <: a \quad q^7 \triangleright \text{ret} <: x1 \quad x1 <: q^7 \triangleright \text{ret}$$

Due to the mutations to `this` and to `x1`, there are constraints in the expected direction and in the opposite direction. Due to line 9, `a` is safe and thus `this` $<$: poly. DFlow creates the following constraints at line 11:

$$a1 <: q^{11} \triangleright \text{this} \quad q^{11} \triangleright \text{this} <: a1 \quad q^{11} \triangleright \text{ret} <: \text{sink1}$$

The last constraint forces `ret` $<$: poly.

If `this` is safe, then due to `a1` $<$: $q^{11} \triangleright \text{this}$, `a1` is forced to be safe, unnecessarily; if `ret` is safe, then due to `x1` $<$: $q^7 \triangleright \text{ret}$, `x1` is forced to be safe, again unnecessarily.

As we argued earlier, `this` being poly is the natural choice: `this.f` is safe in context q^7 and unaffected in q^{11} . Similarly, `ret` being poly is the natural choice: `ret` is safe in context q^{11} and unaffected in q^7 . However, then both q^7 and q^{11} would be safe, and safe would propagate to both `a1` and `x1`.

The problem is that in DFlow, callsite contexts q^i play two roles: (1) they propagate formal-to-return-value dependences to give rise to actual-to-left-hand-side of call assignment dependences, akin to line 14 in CFL-Solver, and (2) they account for N -paths and P -paths, akin to lines 15-18 in CFL-Solver. In most cases, if both the parameter and the return value are poly then there is flow from the parameter to the return value (e.g., `p` $<$: `ret` in `id`), and DFlow is precise. However, when there is no flow from the parameter to the return value, DFlow is imprecise as in Fig. 6.

$$\frac{\text{(TASSIGN)} \quad L \vdash y \dashrightarrow x \quad \text{ReIm}(x) \neq \text{readonly} \Rightarrow L \vdash x \dashrightarrow y \quad \Gamma \vdash L}{\Gamma \vdash x = y}$$

$$\frac{\text{(TWRITE)} \quad L \vdash x \dashrightarrow y.f \quad \text{ReIm}(y.f) \neq \text{readonly} \Rightarrow L \vdash y.f \dashrightarrow x \quad \Gamma \vdash L}{\Gamma \vdash y.f = x}$$

$$\frac{\text{(TREAD)} \quad L \vdash y.f \dashrightarrow x \quad \text{ReIm}(x) \neq \text{readonly} \Rightarrow L \vdash x \dashrightarrow y.f \quad \Gamma \vdash L}{\Gamma \vdash x = y.f}$$

$$\frac{\text{(TCALL)} \quad \text{method}(m) = \text{this}, p \rightarrow \text{ret} \quad L \vdash p \dashrightarrow \text{ret} \Rightarrow L \vdash \text{act}(p) \dashrightarrow x \quad L \vdash p \dashrightarrow p' \wedge \text{ReIm}(p' \text{ at } i) \neq \text{readonly} \Rightarrow L \vdash \text{act}(p) \dashrightarrow \text{act}(p') \quad \text{typeof}(m) = q_{\text{this}}, q_p \rightarrow q_{\text{ret}} \quad \Gamma(x) = q_x \quad \Gamma(y) = q_y \quad \Gamma(z) = q_z \quad q_y <: q_1^i \triangleright q_{\text{this}} \quad \text{ReIm}(\text{this at } i) \neq \text{readonly} \Rightarrow q_1^i \triangleright q_{\text{this}} <: q_y \quad q_z <: q_2^i \triangleright q_p \quad \text{ReIm}(p \text{ at } i) \neq \text{readonly} \Rightarrow q_2^i \triangleright q_p <: q_z \quad q_3^i \triangleright q_{\text{ret}} <: q_x \quad \text{ReIm}(x) \neq \text{readonly} \Rightarrow q_x <: q_3^i \triangleright q_{\text{ret}} \quad \Gamma \vdash L}{\Gamma \vdash i: x = y.m(z)}$$

$$\frac{\text{(LINEARIZE-TREAD)} \quad L \vdash x \dashrightarrow y.f \quad \text{typeof}(f) = \text{poly}}{L \vdash x \dashrightarrow y} \quad \frac{\text{(LINEARIZE-TWRITE)} \quad L \vdash y.f \dashrightarrow x \quad \text{typeof}(f) = \text{poly}}{L \vdash y \dashrightarrow x}$$

$$\frac{\text{(TRANS)} \quad L \vdash x \dashrightarrow y \quad L \vdash y \dashrightarrow z}{L \vdash x \dashrightarrow z} \quad \frac{\text{(WELL-FORMED } L \text{ IN } \Gamma) \quad L \vdash s \dashrightarrow t \Rightarrow \text{typeof}(s) <: \text{typeof}(t)}{\Gamma \vdash L}$$

Figure 7. DFlowCFL type system. L is a system of local flow constraints. We extend `typeof` to retrieve the types of variables and expressions; `typeof(x.f)` = $x \triangleright f$ as expected. (Rule (WELL-FORMED L IN Γ) uses Γ in `typeof`.) `act(p)`, used in (TCALL), returns the actual argument corresponding to parameter p .

5.2 Type System

We now present DFlowCFL, a type system that is equivalent to CFL-Solver and thus strictly more precise than DFlow.

The typing rules for DFlowCFL appear in Fig. 7. DFlowCFL augments DFlow with a set of *local constraints* L and a set of rules for the local constraints in L . Constraints in L are of form $x \dashrightarrow y$, $y.f \dashrightarrow x$ or $x \dashrightarrow y.f$, where x and y are local variables in the *same method* m . $L \vdash x \dashrightarrow y$ means that there is flow from x to y . All constraints in L must obey the subtyping relation, e.g., if $x \dashrightarrow y$ is in L , then $x <: y$ must hold in Γ . This is ensured by rule (WELL-FORMED L IN Γ).

Rule (TRANS) adds constraints $x \dashrightarrow y$ to L due to transitivity. (TRANS) discovers dependences from formal parameters to return values. We call these dependences *method summary constraints*. Rules (LINEARIZE-TREAD) and (LINEARIZE-TWRITE) add constraints $x \dashrightarrow y$ when the adapted field is poly (if f is poly, then $x.f \dashrightarrow y$ entails $x \dashrightarrow y$ because $x \triangleright f$ evaluates to x). These rules “erase” field

```

1: procedure DFlowCFL-INFER
2:   for each variable  $x$  do  $S(x) \leftarrow \{\text{tainted}, \text{poly}, \text{safe}\}$  end for
3:   for each  $q_j^i$  at callsite  $i$  do  $S(q_j^i) \leftarrow \{\text{tainted}, \text{poly}, \text{safe}\}$  end for
4:   for each field  $f$  do  $S(f) \leftarrow \{\text{tainted}, \text{poly}\}$  end for
5:   for each statement  $s$  do Add constraints for  $s$  to  $C$  end for
6:   repeat
7:     for each  $c$  in  $C$  do
8:       SOLVECONSTRAINT( $c$ )
9:       if  $c$  is  $x <: y \triangleright f$  and  $S(f)$  is  $\{\text{poly}\}$  then ▷ Case 1
10:        Add  $x <: y$  into  $C$ 
11:       else if  $c$  is  $x \triangleright f <: y$  and  $S(f)$  is  $\{\text{poly}\}$  then ▷ Case 2
12:        Add  $x <: y$  into  $C$ 
13:       else if  $c$  is  $x <: y$  then ▷ Case 3
14:         for each  $y <: z$  in  $C$  do Add  $x <: z$  to  $C$  end for
15:         for each  $w <: x$  in  $C$  do Add  $w <: y$  to  $C$  end for
16:         for each  $w <: q_j^i \triangleright x$  and  $q_k^i \triangleright y <: z$  in  $C$  do ▷ Case 4
17:           Add  $w <: z$  to  $C$ 
18:         end for
19:       end if
20:     end for
21:   until  $S$  remains unchanged
22: end procedure

```

Figure 8. Type inference for DFlowCFL. Lines 2-5 initialize S and C . Cases 1 and 2 add $x <: y$ into C because $y \triangleright \text{poly}$ always yields y . These cases mirror (LINEARIZE-TREAD) and (LINEARIZE-TWRITE). Case 3 adds constraints due to transitivity. It mirrors (TRANS). Case 4 adds constraints between actual(s) and left-hand-side(s) at calls: if there are constraints $w <: q_j^i \triangleright x$ (flow from actual w to formal x) and $q_k^i \triangleright y <: z$ (flow from return value y to left-hand-side z), and also $x <: y$ (flow from formal to return value), Case 4 adds $w <: z$. Note that line 4 calls SOLVECONSTRAINT(c): the solver infers new constraints, which remove additional infeasible qualifiers from S . This process repeats until S stays unchanged.

information, akin to the way CFL-Solver “erases” field annotations $[f \text{ and }]_f$ when f flows into a sink. Just as with CFL-Solver, erasure happens only if f flows into a sink. Intuitively, constraints $x \dashrightarrow y$ in L correspond to d -paths in G .

DFlowCFL improves over DFlow because its (TCALL) separates role (1), the propagation of formal-to-return-value dependences (e.g., $p \dashrightarrow \text{ret}$), which must give rise to actual-to-left-hand-side dependences ($\text{act}(p) \dashrightarrow \text{act}(\text{ret})$), from role (2), the handling of N -paths and P -paths.

Rule (TCALL) has two parts. The first part propagates all method summary constraints, i.e., formal-to-return-value dependences. For example, in method `id`, there is method summary constraint $p \dashrightarrow \text{ret}$. At call $x = \text{id}(y)$ (TCALL) requires that L entails $y \dashrightarrow x$, as expected. The second part creates constraints that link actuals to formals, such as $q_y <: q_1^i \triangleright q_{\text{this}}$. These constraints capture N -paths and P -paths. Qualifiers q_j^i are distinct, which avoids unnecessary propagation from actuals to left-hand-sides of call assignments.

Let us return to the example in Fig. 6. DFlowCFL enforces the following constraints at line 7:

$a <: q_1^7 \triangleright \text{this}$ $q_1^7 \triangleright \text{this} <: a$ $q_3^7 \triangleright \text{ret} <: x1$ $x1 <: q_3^7 \triangleright \text{ret}$

DFlowCFL enforces the following constraints at line 11:

$a1 <: q_1^{11} \triangleright \text{this}$ $q_1^{11} \triangleright \text{this} <: a1$ $q_3^{11} \triangleright \text{ret} <: \text{sink1}$

Recall that `a` and `sink1` are safe. `this` and `ret` can be poly. $q_1^7 = \text{safe}$, $q_3^7 = \text{tainted}$, $q_1^{11} = \text{tainted}$ and $q_3^{11} = \text{safe}$ satisfy the above constraints. Thus, `a1` and `x1` can be tainted, i.e., unaffected by the safe sink, precisely as in CFL-Solver.

5.3 Type Inference

Given sources and sinks, type inference derives a *valid typing*, i.e. an assignment from program variables to type qualifiers that type checks with the typing rules in Fig. 7.

DFlow and DFlowCFL permit many valid typings. For example, if there are no tainted sources, one can type all variables safe and all fields poly and this assignment type checks. This typing is hardly useful. Intuitively, we would like to infer a typing that *minimizes* the impact of safe sinks through the code, just as CFL-Solver minimizes this impact.

Our type inference algorithm computes the typing that minimizes the impact of safe. It initializes all variables to *sets of qualifiers*. S is the mapping from variables to sets of qualifiers. Programmer-annotated variables, including sources and sinks, are initialized to the singleton set that contains the provided type qualifier. For example, sources and sinks from the annotated library map to $\{\text{tainted}\}$ and $\{\text{safe}\}$, respectively. Fields f are initialized to $S(f) = \{\text{tainted}, \text{poly}\}$. All other variables and all q_j^i are initialized to the maximal set of qualifiers $\{\text{tainted}, \text{poly}, \text{safe}\}$.

The algorithm proceeds to create constraints for all program statements according to the typing rules in Fig. 7. It adds those constraints to set of constraints C . For example, at $x.f = y$ the algorithm adds $y <: x \triangleright f$ to C . If y is not readonly, it adds $x \triangleright f <: y$ as well. Then the algorithm iterates over constraints $c \in C$ and calls SOLVECONSTRAINT(c). SOLVECONSTRAINT(c) removes *infeasible* qualifiers from the sets of variables in c [13]. Consider constraint $c: y <: x$ where $S(y) = \{\text{tainted}\}$ and $S(x) = \{\text{tainted}, \text{poly}, \text{safe}\}$. SOLVECONSTRAINT removes `poly` and `safe` from $S(x)$, because there does not exist a $y \in S(y)$ that satisfies $y <: \text{poly}$ or $y <: \text{safe}$. In the case that the infeasible qualifier is the last element in $S(x)$, the inference terminates with a *type error*.

The type inference algorithm, DFlowCFL-Infer is presented in Fig. 8. It adds local constraints $x <: y$ to C due to transitivity (lines 13-15 mirror rule (TRANS) in Fig. 7) or due to actual-to-left-hand-side “connections” at method calls (lines 16-18 mirror rule (TCALL) in Fig. 7). It keeps removing infeasible qualifiers for each constraint until (1) it terminates with a type error, indicating potential flow from a source to a sink, or (2) it reaches a fixpoint, meaning that there is no flow from sources to sinks. Sect. 7 gives a detailed example.

When the inference terminates without type errors, we derive a valid typing by picking the *maximal* element of $S(x)$ according to the ranking `tainted > poly > safe`. We call this typing the *maximal typing* T . It is a theorem that T *always type checks* with the rules for DFlowCFL in Fig. 7. The preference ranking gives rise to a ranking over all valid typings. We say that $T' > T''$ if T' has a larger number of tainted variables than T'' , or T' and T'' have equal number of tainted variables but T' has a larger number of poly variables than T'' . It is a theorem that T is the *unique maximal typing*⁴ (the proof can be found in [13]). T minimizes the impact of safe through the program, just as CFL-Solver does.

The following theorems relate CFL-Solver to T . Theorem 5.1 states that when T exists, x is inferred as tainted if and only if CFL-Solver discovers no paths from x to sink n .

THEOREM 5.1. *Let T exist.*

1. CFL-SOLVER $\vdash x \xrightarrow{d} n \Leftrightarrow T(x) <: \text{safe}$
2. CFL-SOLVER $\vdash x \xrightarrow{N} n \Leftrightarrow T(x) <: \text{safe}$
3. CFL-SOLVER $\vdash x \xrightarrow{C} n \Leftrightarrow T(x) <: \text{poly}$

⁴ T subsumes every valid typing. Thus, it can be viewed as a principal type.

Theorem 5.2 states that DFlowCFL-Infer terminates with a type error if and only if CFL-Solver discovers a path from a source to a sink. It follows directly from Theorem 5.1.

THEOREM 5.2. $\text{CFL-SOLVER} \vdash x \xrightarrow{C} n$, where x is tainted source $\Leftrightarrow T = \text{error}$.

It is worth noting that T , computed by DFlowCFL-Infer, practically always type checks with DFlow [15, 16]. DFlow is simpler and therefore more appealing than DFlowCFL. This paper argues that even in the rare case when T does not type check with DFlow, there is no flow from sources to sinks.

6. Inheritance and Virtual Calls

The standard approach to handling of virtual calls in CFL-Solver would be to compute a set of actual target methods at each virtual call i , and then run lines 9-11 (Fig. 4) for each target method. One can compute this set by using established techniques such as points-to analysis [18, 23] or CHA [4]. Unfortunately, such handling precludes modularity and compositionality. Suppose we analyzed an open library P using CHA. If we composed P with user code that extended one or more of P 's classes, we would need to reanalyze P because an overriding method could have introduced new flows through callbacks, and changed the analysis result.

We handle inheritance and virtual calls with standard function subtyping augmented with reference immutability. Let method n , where $\text{typeof}(n) = q_{\text{this}_n}, q_{p_n} \rightarrow q_{\text{ret}_n}$ override m , where $\text{typeof}(m) = q_{\text{this}_m}, q_{p_m} \rightarrow q_{\text{ret}_m}$. DFlow requires

$$\begin{array}{ll} q_{\text{this}_m} <: q_{\text{this}_n} & \text{ReIm}(\text{this}_n) \neq \text{readonly} \Rightarrow q_{\text{this}_n} <: q_{\text{this}_m} \\ q_{p_m} <: q_{p_n} & \text{ReIm}(p_n) \neq \text{readonly} \Rightarrow q_{p_n} <: q_{p_m} \\ q_{\text{ret}_n} <: q_{\text{ret}_m} & \text{ReIm}(\text{ret}_m) \neq \text{readonly} \Rightarrow q_{\text{ret}_m} <: q_{\text{ret}_n} \end{array}$$

The above constraints guarantee that if there is a pair of parameters of overriding method n such that $p_n <: p'_n$, then the same subtyping holds for the corresponding pair of parameters of overridden method m : $p_m <: p'_m$. (Recall that p' is usually ret .) For example, suppose that n has statement return p but m does not. This entails $p_n <: \text{ret}_n$. Due to the above constraints, we have $p_m <: p_n$ (contravariant arguments) and $\text{ret}_n <: \text{ret}_m$ (covariant return values), which yields $p_m <: p_{\text{ret}}$. Thus, return p in overriding method n will give rise to the expected actual-to-left-hand-side dependence at calls to overridden method m . This handling is modular, in the sense that we can analyze any given set of classes P . We can compose user code with P without reanalysis of P — we simply must check that the above constraints hold for every n and m , where n is a user method that overrides m from P .

DFlowCFL imposes the following requirement, in addition to the above subtyping constraints:

$$L \vdash p_n \dashrightarrow p'_n \Rightarrow L \vdash p_m \dashrightarrow p'_m$$

DFlowCFL requires explicit method summary constraints and the above implication ensures that calls to m give rise to actual-to-left-hand-side dependences triggered by n . The subtyping constraints capture N -paths and P -paths.

Finally, CFL-Solver calls $\text{BIEDGE}(\text{this}_n, \text{this}_m, \text{this}_n, d, d)$, $\text{BIEDGE}(p_n, p_m, p_n, d, d)$ and $\text{BIEDGE}(\text{ret}_m, \text{ret}_n, \text{ret}_m, d, d)$. Thus, CFL-Solver creates edges $\text{this}_m \rightarrow \text{this}_n$, from $p_m \rightarrow p_n$ and from $\text{ret}_n \rightarrow \text{ret}_m$ plus the corresponding inverse edges when necessary.

7. Example with Taint Analysis for Android

In this section, we demonstrate an example of the workings of DFlowCFL-Infer for the detection of privacy leaks in Android apps. Consider the FieldSensitivity2 example in Fig. 9. This example is refactored from DroidBench [10], a set of micro benchmarks for

taint analysis for Android; it is representative of real-world leaks. The return of `TelephonyManager.getSimSerialNumber` (line 10) is a *source* and the parameter `msg` of `SmsManager.sendMessage` (line 16) is a *sink*. The serial number of the SIM card is obtained and stored into a `Data` object. Later, it is retrieved from the `Data` object and sent out through an SMS message without user consent. Fig. 9 demonstrates DFlowCFL-Infer (Fig. 8).

The interpretation of DFlow in terms of CFL-reachability is very important for error reporting. As mentioned in Sect. 2, each type error can be explained with one or more CFL-reachability source-sink flow paths. Fig. 9 presents the flow path that is to blame for the type error in Fig. 9.

8. Related Work

Shankar et al. present an integrity type system for detecting string format vulnerabilities in C programs [31]. The type system has two type qualifiers, tainted and untainted; polymorphism is not part of the core system. In contrast, DFlow and DFlowCFL-Infer handle polymorphism naturally. Ernst et al. present IFC, an integrity type system for detecting privacy leaks in Android apps [8]. IFC is similar to DFlow but it requires annotations and works only on Java source. In contrast to previous work on integrity types, this paper focuses on the connection of integrity types to CFL-reachability.

Rehof and Fahndrich [26] study the connection between type-based flow analysis with polymorphic subtyping and CFL-reachability. They do not allow for polymorphism in the type structure, which essentially entails graphs without brackets. In contrast, we handle graphs with both parentheses and brackets. Furthermore, Rehof and Fahndrich [26] do not discuss mutable references. Fahndrich et al. [9] present an application of [26] to context-sensitive Steensgard-style points-to analysis for C. This work handles mutable references, because it uses equality (i.e., unification) constraints. Equality constraints is the standard approach to the handling of mutable references [11, 29, 31]. However, this approach is imprecise [22]. A key novelty in our work is the use of reference immutability to allow for limited subtyping.

Sridharan and Bodik [32] present refinement-based points-to analysis for Java using CFL-reachability. Recently, Shang et al. [30] and Lu et al. [21] build upon this work towards an incremental demand-driven points-to analysis for Java. Xu et al. [35] improve the scalability of CFL-reachability-based points-to analysis. These works focus on points-to analysis and all require heap abstraction. Heap abstraction requires handling of reflective object creation, which is inherently difficult; also, it may preclude modularity and compositionality. One key achievement of our work is the CFL-reachability-based integrity analysis, which avoids heap abstraction and is modular and compositional. Another key difference with previous work is our focus on type-based approaches and their connection to CFL-reachability.

There is a lot of recent work on static taint analysis for Android [10, 12, 17, 20, 36]. These analyses typically use points-to and dataflow-based approaches. Type-based taint analysis appears to be better suited to the problem.

Volpano et al. [34] and Myers [24] present type systems for secure information flow. These systems are substantially more complex and powerful than DFlow. They focus on type checking and do not include type inference, or include only local type inference.

We conclude with a concrete discussion of this work in relation to our previous work in [15] and [16]. The work in [15], our first effort on type-based integrity analysis, presents the SFlow type system, the corresponding type inference analysis SFlowInfer, and their application to taint analysis for Java web applications. SFlow is similar to DFlow; however SFlow is strictly less precise than DFlow. [15] presents an empirical evaluation of SFlowInfer over a wide range of Java web applications. The work in [16] builds upon [15].

```

1 public class Data {
2   String secret;
3   String get(Data this) {return this.secret;}
4   void set(Data this, String p){this.secret = p;}
5 }
6 public class FieldSensitivity2 extends Activity {
7   protected void onCreate(Bundle b) {
8     Data dt = new Data();
9     TelephonyManager tm = (TelephonyManager)
10      getSystemService("phone");
11     String sim = tm.getSimSerialNumber();
12     dt.set(sim);
13     SmsManager sms = SmsManager.getDefault();
14     String sg = dt.get();
15
16     sms.sendMessage("+123", null, sg, null, null);
17   }
18 }

```

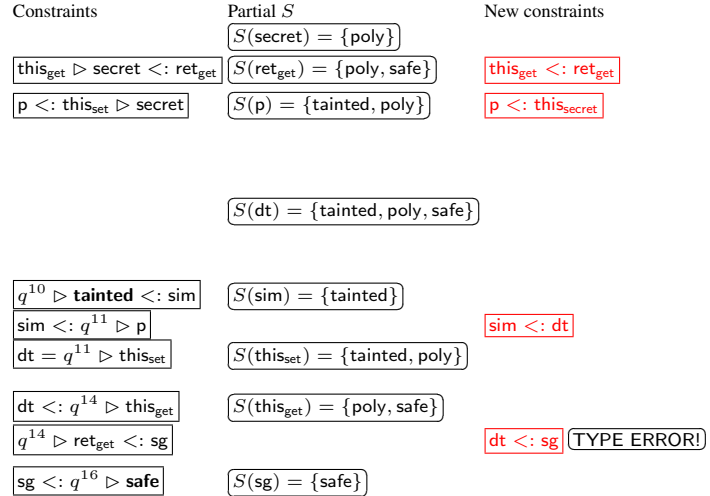


Figure 9. FieldSensitivity2 example refactored from DroidBench. The frame box beside each statement shows the corresponding constraints the statement generates. We omitted uninteresting constraints. The oval boxes show initial propagation. 16 forces sg to be $\{\text{safe}\}$, then 14 forces ret_{get} to be $\{\text{poly, safe}\}$ and then 3 forces this_{get} to be $\{\text{poly, safe}\}$ and secret to be $\{\text{poly}\}$. 10 forces sim to be $\{\text{tainted}\}$, which in turn forces the parameters p and this_{set} to be $\{\text{tainted, poly}\}$. The red frame box in the fourth column (New constraints) shows the added constraints (due to lines 9-19 in Fig. 8). Since field secret is poly , constraint $\text{this}_{\text{get}} \triangleright \text{secret} <: \text{ret}_{\text{get}}$ leads to $\text{this}_{\text{get}} <: \text{ret}_{\text{get}}$ (Case 2), which in turn leads to $\text{dt} <: \text{sg}$ due to the call at 14 (Case 4). Similarly, $p <: \text{this}_{\text{set}} \triangleright \text{secret}$ leads to $p <: \text{this}_{\text{secret}}$, which in turn leads to $\text{sim} <: \text{dt}$ due to the call at 11. Since sim is $\{\text{tainted}\}$ and sg is $\{\text{safe}\}$, these constraints cause a TYPE ERROR, detecting the leak. We envision that the type error will be reported at the source or at the sink and will come with this source-sink path: $\text{sim} \xrightarrow{(11)} p \xrightarrow{\{\text{secret}\}} \text{this}_{\text{set}} \xrightarrow{(14)} \text{dt} \xrightarrow{\{\text{poly, safe}\}} \text{this}_{\text{get}} \xrightarrow{\{\text{poly, safe}\}} \text{ret}_{\text{get}} \xrightarrow{(16)} \text{sg}$.

It extends [15] in three directions. First, it develops DFlow, which is more precise than SFlow. Second, it presents novel handling of Android-specific features. Third, the tools presented in [16] work on APKs, while SFlowInfer works only on Java source; thus, we analyze 73 popular apps from the Google Play Store and detect interesting leaks. This previous work, especially [16], motivated the current paper. Both [15] and [16] focus on the respective type systems and extensive experimentation with those type systems and the corresponding type inference tools. In the paper, we establish a connection between the type-based approaches from [15] and [16] and CFL-reachability, which is a well-established precise program flow analysis technique. The connection to CFL-reachability helps error reporting, as we argued earlier. In addition, it helps explain the approximations employed by the type-based approaches, in terms of Reprs' framework [28].

9. Conclusion

We presented DFlow, a context-sensitive integrity type system and we give an interpretation of DFlow in terms of CFL-reachability. We proposed DFlowCFL, a more precise integrity type system, and DFlowCFL-Infer, the corresponding type inference analysis, which is equivalent to CFL-reachability. DFlowCFL-Infer is an effective taint analysis for Android detecting numerous privacy leaks in Android apps from the Google Play Store.

Acknowledgements We thank the anonymous reviewers for their helpful feedback. This work was supported by NSF Award CCF-1319384 and a Google Faculty Research Award (February 2013).

References

[1] Matt Welsh, Volatile and Decentralized. http://matt-welsh.blogspot.com/2013_01_01_archive.html, 2013.

- [2] Tainting Checker. <http://types.cs.washington.edu/checker-framework/current/checkers-manual.html#tainting-checker>, 2014.
- [3] S. Arzt, S. Rasthofer, C. Fritz, E. Bodden, A. Bartel, J. Klein, Y. le Traon, D. Octeau, and P. McDaniel. FlowDroid: Precise context-, flow-, field-, object-sensitive and lifecycle-aware taint analysis for android apps. In *PLDI, to appear*, 2014.
- [4] J. Dean, D. Grove, and C. Chambers. Optimization of object-oriented programs using static class hierarchy analysis. In *ECOOP*, pages 77–101, 1995.
- [5] W. Dietl and P. Müller. Universes: Lightweight ownership for JML. *Journal of Object Technology*, 4(8):5–32, 2005.
- [6] W. Enck, P. Gilbert, B.-G. Chun, L. P. Cox, J. Jung, P. McDaniel, and A. N. Sheth. TaintDroid: An information-flow tracking system for realtime privacy monitoring on smartphones. In *OSDI*, pages 1–6, 2010.
- [7] M. D. Ernst. Type Annotations specification (JSR 308). <http://types.cs.washington.edu/jsr308/>, 2012.
- [8] M. D. Ernst, R. Just, S. Millstein, W. M. Dietl, S. Pernsteiner, F. Roesner, K. Koscher, P. Barros, R. Bhoraskar, S. Han, P. Vines, and E. X. Wu. Collaborative verification of information flow for a high-assurance app store. Technical Report UW-CSE-14-04-02, University of Washington Department of Computer Science and Engineering, Apr. 2014.
- [9] M. Fähndrich, J. Rehof, and M. Das. Scalable context-sensitive flow analysis using instantiation constraints. In *PLDI*, pages 253–263, 2000.
- [10] C. Fritz, S. Arzt, S. Rasthofer, E. Bodden, J. Klein, Y. le Traon, D. Octeau, and P. McDaniel. Highly precise taint analysis for Android application. Technical Report TUD-CS-2013-0113, EC SPRIDE, 2013.
- [11] R. M. Fuhrer, F. Tip, A. Kiezun, J. Dolby, and M. Keller. Efficiently refactoring java applications to use generic libraries. In *ECOOP*, pages 71–96, 2005.
- [12] C. Gibler, J. Crussell, J. Erickson, and H. Chen. AndroidLeaks: Automatically detecting potential privacy leaks in Android applications on a large scale. In *TRUST*, pages 273–290, 2012.

- [13] W. Huang, W. Dietl, A. Milanova, and M. D. Ernst. Inference and checking of object ownership. In *ECOOP*, pages 181–206, 2012.
- [14] W. Huang, A. Milanova, W. Dietl, and M. D. Ernst. ReIm & ReImInfer: Checking and inference of reference immutability and method purity. In *OOPSLA*, pages 879–896, 2012.
- [15] W. Huang, Y. Dong, and A. Milanova. Type-based taint analysis for Java web applications. In *FASE*, pages 140–154, 2014.
- [16] W. Huang, Y. Dong, and A. Milanova. DFlow & DroidInfer: Effective Type-based Taint Analysis for Android. Submitted for publication, 2014. URL <http://www.cs.rpi.edu/~milanova/docs/DroidInfer.pdf>.
- [17] J. Kim, Y. Yoon, and K. Yi. SCANDAL: Static analyzer for detecting privacy leaks in Android applications. In *MoST*, 2012.
- [18] O. Lhoták and L. Hendren. Scaling Java points-to analysis using Spark. In *CC*, pages 153–169, 2003.
- [19] S. Liang, A. W. Keep, M. Might, S. Lyde, T. Gilray, and P. Aldous. Sound and precise malware analysis for Android via pushdown reachability and entry-point saturation. In *SPSM*, pages 21–32, 2013.
- [20] L. Lu, Z. Li, Z. Wu, W. Lee, and G. Jiang. CHEX: Statically vetting Android apps for component hijacking vulnerabilities. In *CCS*, pages 229–240, 2012.
- [21] Y. L. Lu, L. Shang, X. Xie, and J. Xue. An incremental points-to analysis with cfl-reachability. In *CC*, pages 61–81, 2013.
- [22] A. Milanova and W. Huang. Composing polymorphic information flow systems with reference immutability. In *FTJJP*, pages 5:1–5:7, 2013.
- [23] A. Milanova, A. Rountev, and B. G. Ryder. Parameterized object sensitivity for points-to analysis for Java. *ACM Transactions on Software Engineering and Methodology*, 14(1):1–41, Jan. 2005.
- [24] A. C. Myers. JFlow: Practical mostly-static information flow control. In *POPL*, pages 228–241, 1999.
- [25] A. C. Myers, J. A. Bank, and B. Liskov. Parameterized types for Java. In *POPL*, pages 132–145, 1997.
- [26] J. Rehof and M. Fähndrich. Type-based flow analysis: From polymorphic subtyping to CFL reachability. In *POPL*, pages 54–66, 2001.
- [27] T. Reps. Program analysis via graph reachability. *Information and Software Technology*, 40:5–19, 1998.
- [28] T. Reps. Undecidability of context-sensitive data-independence analysis. *ACM Transactions on Programming Languages and Systems*, 22(1):162–186, 2000.
- [29] A. Sampson, W. Dietl, and E. Fortuna. EnerJ: Approximate data types for safe and general low-power computation. In *PLDI*, pages 164–174, 2011.
- [30] L. Shang, Y. Lu, and J. Xue. Fast and precise points-to analysis with incremental cfl-reachability summarisation: preliminary experience. In *ASE*, pages 270–273, 2012.
- [31] U. Shankar, K. Talwar, J. S. Foster, and D. Wagner. Detecting format string vulnerabilities with type qualifiers. In *USENIX Security*, pages 201–220, 2001.
- [32] M. Sridharan and R. Bodík. Refinement-based context-sensitive points-to analysis for java. In *PLDI*, pages 387–400, 2006.
- [33] M. S. Tschantz and M. D. Ernst. Javari: Adding reference immutability to Java. In *OOPSLA*, pages 211–230, 2005.
- [34] D. Volpano, G. Smith, and C. Irvine. A sound type system for secure flow analysis. *Journal of Computer Security*, pages 167–187, 1996.
- [35] G. H. Xu, A. Rountev, and M. Sridharan. Scaling cfl-reachability-based points-to analysis using context-sensitive must-not-alias analysis. In *ECOOP*, pages 98–122, 2009.
- [36] Z. Yang and M. Yang. LeakMiner: Detect information leakage on Android with static taint analysis. *2012 Third World Congress on Software Engineering*, pages 101–104, Nov. 2012.