

# Implementation and Evaluation of a Static Backwards Data Flow Analysis in FlowDroid

Implementierung und Evaluation einer statischen rückwärtsgerichteten Datenflussanalyse in FlowDroid

Bachelor thesis by Tim Lange

Date of submission: 22nd March 2021

1. Review: Dr. Steven Arzt

2. Review: Prof. Dr. Michael Waidner  
Darmstadt



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT



**Fraunhofer**  
SIT

Computer Science  
Department  
Fraunhofer SIT  
Secure Software  
Engineering

---

# Contents

---

<b>1. Introduction</b>	<b>5</b>
<b>2. Background</b>	<b>6</b>
2.1. Static Data Flow Analysis . . . . .	6
2.2. IFDS . . . . .	7
2.2.1. Original Definition . . . . .	7
2.2.2. Practical Extensions . . . . .	10
2.3. Access Paths . . . . .	11
2.4. Intermediate Representations . . . . .	12
2.5. FlowDroid . . . . .	13
<b>3. Theory</b>	<b>15</b>
3.1. Flow Functions . . . . .	15
3.1.1. Normal Flow . . . . .	15
3.1.2. Call Flow . . . . .	16
3.1.3. Return Flow . . . . .	18
3.1.4. CallToReturn Flow . . . . .	19
3.2. Complexity of Data Flow Analysis . . . . .	20
<b>4. Implementation</b>	<b>23</b>
4.1. Integration . . . . .	23
4.2. Flow-Sensitive Alias Analysis . . . . .	24
4.3. Rules . . . . .	25
4.3.1. Source & Sink Propagation Rule . . . . .	26
4.3.2. Backwards Array Propagation Rule . . . . .	26
4.3.3. Backwards Exception Propagation Rule . . . . .	26
4.3.4. Backwards Wrapper Propagation Rule . . . . .	27
4.3.5. Backwards Strong Update Rule . . . . .	27
4.3.6. Backwards Clinit Rule . . . . .	29

---

4.3.7. Other Rules . . . . .	30
4.4. Other Components . . . . .	30
4.4.1. Code Optimizer: AddNOPStmts . . . . .	30
4.4.2. Taint Wrappers . . . . .	31
4.4.3. Native Call Handler . . . . .	31
<b>5. Validation</b>	<b>33</b>
5.1. Unit Tests . . . . .	33
5.2. DroidBench . . . . .	34
5.2.1. Configuration . . . . .	34
5.2.2. Results . . . . .	35
5.2.3. Results Explanation . . . . .	39
5.2.4. Improvements From The Summary Taint Wrapper . . . . .	41
<b>6. Performance Evaluation</b>	<b>43</b>
6.1. DroidBench . . . . .	43
6.1.1. Results . . . . .	44
6.1.2. Result Explanation . . . . .	47
6.1.3. Using A More Precise Taint Wrapper . . . . .	48
6.2. Real World Apps . . . . .	52
6.2.1. Configuration . . . . .	52
6.2.2. Results . . . . .	53
6.2.3. Comparison to forwards analysis . . . . .	53
<b>7. Related Work</b>	<b>58</b>
<b>8. Conclusion</b>	<b>60</b>
<b>Bibliography</b>	<b>61</b>
<b>A. Appendix</b>	<b>63</b>

---

## **Erklärung zur Abschlussarbeit gemäß §22 Abs. 7 APB TU Darmstadt**

---

Hiermit versichere ich, Tim Lange, die vorliegende Bachelorarbeit gemäß §22 Abs. 7 APB der TU Darmstadt ohne Hilfe Dritter und nur mit den angegebenen Quellen und Hilfsmitteln angefertigt zu haben. Alle Stellen, die Quellen entnommen wurden, sind als solche kenntlich gemacht worden. Diese Arbeit hat in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde vorgelegen.

Mir ist bekannt, dass im Falle eines Plagiats (§38 Abs. 2 APB) ein Täuschungsversuch vorliegt, der dazu führt, dass die Arbeit mit 5,0 bewertet und damit ein Prüfungsversuch verbraucht wird. Abschlussarbeiten dürfen nur einmal wiederholt werden.

Bei einer Thesis des Fachbereichs Architektur entspricht die eingereichte elektronische Fassung dem vorgestellten Modell und den vorgelegten Plänen.

Darmstadt, 22nd March 2021

---

T. Lange



---

# 1. Introduction

---

---

## 2. Background

---

In this chapter, we introduce the necessary background. In section 2.1, we explain the term static data flow analysis. We introduce concepts used to solve data flow problems precisely in subsection 2.2.1 and section 2.3. We reason the need for a more manageable code representation in section 2.4 and last, we introduce FLOWDROID, the tool our work is based on, in section 2.5.

---

### 2.1. Static Data Flow Analysis

---

In the field of compilers, there is a distinction between static and dynamic. Static generally refers to something that is decided at compile-time, while dynamic refers to runtime decisions [1]. The same distinction is also present for analyses. Dynamic analysis observes the program's run behavior, while static analysis works on a representation of the code. Both have different tradeoffs. Achieving good code coverage is hard in dynamic analysis and it is hard only to traverse the actually taken paths in static analysis [2]. Thus dynamic analysis is an underapproximation and static analysis is an overapproximation. In the following, we only consider static analysis.

Data flow analysis is a broad term for analyses that try to identify data flows through a program. Khedker [9] defines data flow analysis as follows:

*Data flow analysis is a process of deriving information about the run time behavior of a program.*

Data flow analyses are used in many different ways. Compilers use it to apply optimizations, others use it for software verification and it is also used for reverse-engineering [9]. A special kind of data flow analysis is taint analysis which falls under the category of reverse-engineering and its concepts might be familiar from code reviews. In taint analysis, the goal is to determine whether a particular statement's information contents flow through

---

the program to another statement. Variables that contain such valuable information are tainted variables. This valuable information has to come from somewhere, the so-called sources. Sources can be any expression but are often methods and their returned values are considered tainted. On the other end, sinks leak valuable information. A data flow between a source and a sink is called a leak [2]. If taint analysis is applied to the use case of tracking detection, a source would be a method returning a unique identifier and a sink would be a method that sends out data to the internet. When finding a leak, we know the receiving server can identify the device.

There is also a categorization for data flow analyses. Sensitivities describe whether an analysis is capable of considering an aspect. There are five common sensitivities:

- **Flow Sensitivity:** A flow-sensitive analysis can determine if a fact holds at a particular statement. Data flow analyses are always flow-sensitive [9].
- **Context Sensitivity:** An interprocedural analysis can distinguish the context of a called method, e.g., knows the original call site at a return statement.
- **Object Sensitivity:** An analysis can distinguish field accesses on different objects.
- **Field Sensitivity:** An analysis can distinguish different field accesses on the same object.
- **Path Sensitivity:** An analysis takes conditional branches into account, e.g., the condition holds after the branch.

Is this common knowledge? Didn't find a good source for this

We also need a representation for the information the analysis gathered: the data flow fact. A data flow fact is a logical assertion that is either true or false at a statement. Now, there are two different kinds of facts: may and must. For a must analysis, the fact must hold on all paths to this statement, while a may analysis only guarantees the fact holds on one path. The decision of which kind fits depends on the type of data flow analysis. Taint analyses like FLOWDROID are based on the may analysis [2].

---

## 2.2. IFDS

---

### 2.2.1. Original Definition

Interprocedural finite distributive subset (IFDS) problems are a special class of a data flow analysis problem. Generally, the solution to a data flow problem is the meet-over-

---

all-paths (MOP) solution, which is undecidable [16]. But all problems adhering to IFDS can be transformed into a graph-reachability problem and consequently, the solution is computable in polynomial time. It is context-sensitive and flow-sensitive by default [15].

IFDS operates on a so-called exploded supergraph. Every node in the exploded supergraph is a tuple  $\langle s, d \rangle$  of a statement  $s$  in the interprocedural control-flow graph and a data flow fact  $d$ . The domain is typically the set of variables in the program. Edges between two nodes  $\langle s, d \rangle$  and  $\langle s', d' \rangle$  exist if  $d$  propagated over  $s$  yields  $d'$  and  $s'$  is a successor of  $s$ . Propagating facts along the control-flow graph already ensures flow-sensitivity. For context-sensitivity, IFDS only visits valid paths. Reps et al. proposed a context-sensitive grammar that acts as a call stack to ensure no mismatch and the path is a valid execution path [15].

To propagate facts over statements, we need to define rules on how the data flow changes when observing a statement. These rules are called flow functions. There are four types of flow functions:

- **Call Flow:** Edges from call statement into a method. The flow function maps the facts visible in the callee into it [15].
- **Return Flow:** Edges returning from a method. The flow function maps the facts visible in the caller out of the callee [15].
- **Call To Return Flow:** Edges over a call statement. The flow function maps the facts not visible in the callee over the call statement [15].
- **Normal Flow:** Edges over every other statement. Often, this flow function only handles assign statements [15].

The incoming set of facts is all predecessors' outgoing facts merged using a merge operator  $\sqcap$ :

$$in(s) := \bigsqcap_{p \in Preds(s)} out(p)$$

Now, we also want to introduce new facts. For that reason, the domain contains a zero fact and all nodes with  $d = \mathbf{0}$  are always reachable; thus, the zero fact is a tautology. Whenever we want to introduce a fact, we can model this in the flow function by deriving such facts from the zero fact [15]. For example, in taint analysis, the flow functions map zero facts at sources to a tainted variable.



---

IFDS also utilizes summaries. After returning from a method, the algorithm solved a subproblem for which it remembers the results to be applied later. So, the proposed tabulation algorithm for solving the realizable path problem is a dynamic algorithm [15].

Eventually, there is no fact to propagate anymore and the analysis will terminate. Facts die because a flow function killed them or the algorithm already observed the same fact at a statement and reached a fixpoint [15].

However, we already started this section, hinting not all problems can be formulated in the IFDS framework. The restrictions the problems have to abide by are eponymous in IFDS and explained in the following paragraphs.

**Distributive** The flow function must be distributive over the merge operator. Formally,  $f(x \sqcap y) = f(x) \sqcap f(y)$  must hold at any time. Informally speaking, it does not matter whether facts get merged before or after applying the flow functions. By defining the flow function signature as  $f : Fact \rightarrow Facts$  with a single fact as input but a set of facts as output, this property is trivially satisfied. Distributiveness is essential for the correctness of IFDS, because only if the flow functions are distributive the maximum fixed point (MFP) equals MOP and MFP is computable in polynomial time [9, 15].

**Finite** Another restriction is that the set of data flow facts has to be finite. Let us go by a counterexample of what IFDS is not capable of: Answering "Which value is stored in variable  $x$  at statement  $s$ ?". Now the data flow fact is a tuple of the variable together with the stored value  $\langle x, v \rangle$ . Consider Figure 2.1. Assume  $x$  is an integer of infinite precision for the domain to be infinite.  $x$  is initialized to zero and passed into the method `foo()` multiple times. Now consider the summaries in Figure 2.1b. We never get to use one summary because the value always increases. Because of the same reason, we also never reach a fixpoint. Thus, the domain has to be finite and, in practice, also small as the domain is cubic in the time-complexity  $O(|E| \cdot |D|^3)$  [15].

**Subset** Data flow frameworks need to deal with merging the outcoming sets to a single incoming set. Essentially, to formalize the approximation and satisfy ordering constraints, data flow frameworks rely on lattices [9]. IFDS also defines an underlying lattice on the powerset of the domain. The lattice ordering must be set inclusion. Therefore, the merge operator is set union or set intersect. Now recall may and must from the last subsection. Here we can see the connection between the merge operator and may or must.

---

---

1	<code>void main() {</code>	
2	<code>    RealInteger x = 0;</code>	
3	<code>    while (condition)</code>	
4	<code>        x = foo(x);</code>	$\langle x, 0 \rangle \rightarrow \langle x, 42 \rangle$
5	<code>}</code>	$\langle x, 42 \rangle \rightarrow \langle x, 84 \rangle$
6		$\langle x, 84 \rangle \rightarrow \langle x, 126 \rangle$
7	<code>RealInteger foo(RealInteger x) {</code>	$\dots$
8	<code>    return x + 42;</code>	
9	<code>}</code>	
	(a) Code	(b) Summaries

Figure 2.1.: Finiteness example

The paper by Reps et al. later decides on set union due to the duality of must and may not [15]. This decision is also efficient in practice, as discussed in the following subsection.

### 2.2.2. Practical Extensions

The original definition is inefficient in practice. Among others, Naeem et al. proposed practical extensions to the IFDS framework to perform better in practice [13].

The original algorithm demands a fully built exploded supergraph. Even in moderate programs, the domain can get quite large. As the nodes in the exploded supergraph are the cross-product of the domain and interprocedural call-graph nodes, it is infeasible to generate the full graph beforehand. Because there is no way to know before which part of the supergraph the analysis demands, they propose to generate it ad-hoc. That also removes the restriction on a small domain. Now IFDS is also feasible if the domain's encountered subset is small enough [13]. The restrictions on the domain set can be loosened even more. Bodden suggests in-practice, the domain can be infinite. Only the observed facts must adhere to the ascending-chain condition over the flow functions when using the on-demand supergraph [4].

Also, the original IFDS definition ignores the type structure of the programming language. It can be used to kill facts due to impossible casts. Also, facts with the same variable but different types can be merged with the superclass as its new type [13].

---

Furthermore, the original definition starts the IFDS algorithm at the entry point of the interprocedural call-graph. As described in subsection 2.2.1, a flow function can derive an initial fact from the zero fact whenever needed. If the methods where initial facts will be introduced are known a priori, the supergraph can be traversed without applying flow functions until such a method is found on the path. This optimization introduces unbalanced problems where a method return but no corresponding call site is found, which can be solved by a small extension to the tabulation algorithm. The extension was first described by Lerch [12] and is also present in FLOWDROID [2].

Another optimization is possible if the merge operator is set union. Then there is no need to wait for other predecessors to finish as a set is always a subset of a union with itself and another set ( $A \subseteq A \cup B$  holds). Hence, the IFDS solver can skip the *in*-set construction and immediately propagate the outgoing facts as singleton sets, which is beneficial in a parallelized solver [17].

Is there a name for this property?

---

## 2.3. Access Paths

---

We have already seen IFDS fulfills context- and flow-sensitivity by default. Now, a precise analysis for Java also needs object- and field-sensitivity. Thus, we also need to model the heap.

Access paths are one possible heap model. They consist of a list of field dereferences linked to a tainted variable of a reference type [9]. Note, this increases the domain size because now not only "object *o* is tainted" is a data flow fact, but also all of its fields can be tainted. Especially when encountering recursive data structures such as linked lists, this gets problematical. Consider Figure 2.2, the loop would let the observed domain grow indefinitely and never reach a fixpoint. As a solution, access paths are limited in length, which is also called *k*-limiting, whereas the constant *k* is the maximum access path length. If an access path passes this length, it is cut off and the entire last reference is considered tainted [8]. Consider again Figure 2.2 with  $k = 2$ . This time the analysis would reach the fixpoint *lst.next.prev.\** after two iterations. The cut-off comes with a loss of precision but is inevitable.

Although with *k*-limiting, the algorithm terminates again, it introduces a new problem. After a loop like in Figure 2.2, the access path is polluted with a dereference chain to its base object even though the *next.prev* dereference could be omitted without precision loss. Thus, Deutsch proposed symbolic access paths, which try to eliminate loops in access

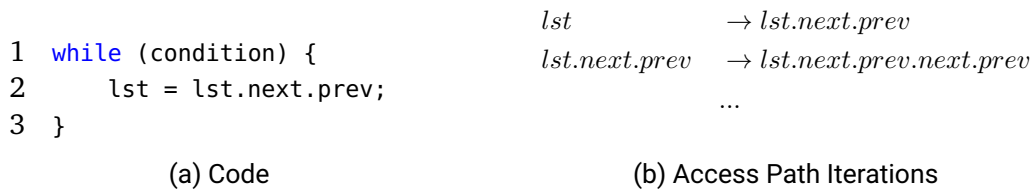


Figure 2.2.: Infinite Access Path

paths [6]. In practice, Deutsch’s approach needs some adaptations as he only considered fields but not base objects and he defines loops simply by type [2].

---

## 2.4. Intermediate Representations

---

Most compilers these days use intermediate representations (IRs). IRs are an equivalent representation of the source code but are much simpler and more regular and are typically not architecture-dependent. They are often in an interchangeable format and can be saved as text to be used by various tools [19]. Such an IR allows compilers to apply machine-independent optimizations to the code without worrying about complex expressions in the source code or reimplementing the optimization for each architecture.

The Java Virtual Machine (JVM) also operates on an IR called Java bytecode. The JVM is mostly stack-based and so is the Java bytecode. In Figure 2.3 is an example of a simple code snippet translated to Java bytecode. Simple expressions such as `c = a + b` translate into multiple statements and there is no fixed length of an expression in the bytecode. The analysis would also have to reconstruct the expressions ad-hoc. Furthermore, Java bytecode has over 200 possible instructions<sup>1</sup>, which need to be considered and only knows primitive types and references. Concluding, stack-based IRs are suitable for just-in-time interpretation but inconvenient for data flow analysis [20].

A more convenient representation for static analysis is three-address codes. Each statement consists of up to three operands and is either an assignment or a control-flow statement. Such a representation is closer to the original source code while reducing the possible combinations to a manageable amount [1].

---

<sup>1</sup><https://docs.oracle.com/javase/specs/jvms/se8/html/>

---

1	<code>int a = 21;</code>	1	<code>bipush 21 // push 21</code>
2	<code>int b = 21;</code>	2	<code>istore_1 // store in register 1</code>
3	<code>int c = a + b;</code>	3	<code>bipush 21 // push 21</code>
		4	<code>istore_2 // store in register 2</code>
		5	<code>iload_1 // push a</code>
		6	<code>iload_2 // push b</code>
		7	<code>iadd // pop a &amp; b and push a + b</code>
		8	<code>istore_3 // store in register 3</code>
	(a) Java code		(b) Java bytecode

Figure 2.3.: Java bytecode example

Jimple is a three-address intermediate representation and can be constructed from the Java and Dalvik bytecode, the IR used for Android apps. It is a high-level representation and its syntax is close to Java. Complex statements are split up into multiple statements. For example, there can be only one field reference per statement and arguments are always local variables. Jimple also reconstructs reference types [20]. This groundwork greatly reduces the possible cases the data flow analysis needs to consider and eases the analysis.

---

## 2.5. FlowDroid

---

FLOWDROID is a precise context-, flow-, object- and field-sensitive static taint analysis tool for Android apps[3]. Since its initial release in 2014, it is actively maintained and gained traction in research and academia<sup>2</sup>. It is based on Soot, a Java optimization framework, which later has been extended for static analysis [10]. Soot provides the call graph and the conversion from Java and Dalvik bytecode to Jimple, the intermediate representation of choice for FLOWDROID [2].

Androids activity-lifecycle concept does not have a single entry point; instead, there are multiple callbacks as a possible entry point. Also, an Android app can contain multiple components and register callbacks in various of Android's standard libraries. FLOWDROID

---

<sup>2</sup><https://github.com/secure-software-engineering/FlowDroid>

models the entire Android lifecycle to be precise and generates a dummy main method to provide a single entry point for the call graph generation. [3].

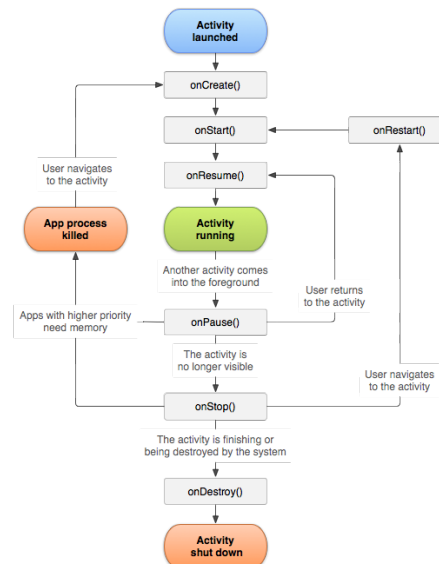


Figure 2.4.: Activity Lifecycle<sup>3</sup>

To provide precise results even with aliases in use, FLOWDROID contains, besides the taint analysis, another IFDS problem to resolve all encountered aliases on-demand. For this, FLOWDROID needs a high object-sensitivity. Hence it makes use of symbolic access paths [3] .

The implementation of FLOWDROID is modular, easily extensible and offers many additional features. Two of them are noteworthy for this work: native call handler and taint wrappers. As both Java and Android allow calling native methods, FLOWDROID also needs to model those cases. It currently does not support the analysis of those methods but contains rules for essential methods. The second feature is taint wrappers. They allow defining rules for methods, e.g., from a commonly used feature such as StringBuilder, which allows the taint analysis to skip the method and apply a summary [3]. STUBDROID, an extension to FlowDroid by Arzt et al., allows precomputing summaries using FLOWDROID and serializes them in an XML format for tool-independent use. These summaries are handy for real-world applications where often third-party libraries are used [2].

is sym-  
bolic  
access  
paths  
right?

<sup>3</sup>Taken from <https://developer.android.com/guide/components/activities/activity-lifecycle>.

---

## 3. Theory

---

In the first part of this chapter, we define the flow functions and in the second part we discuss the time complexity of taint analysis specifically highlighting differences between forward and backward analysis.

---

### 3.1. Flow Functions

---

We describe the behavior of the flow functions based on the Jimple language and define semi-formal rules analogous to the publication[2] on FLOWDROID. These rules only focus on the base control-flow, we describe flows for additional language features such as arrays and exceptions later more informal in section 4.3. Note that traversing the interprocedural control-flow graph backwards means call edges are now return edges and vice versa.

#### 3.1.1. Normal Flow

Normal flow functions handle every statement that does not contain an `InvokeExpr` and assignments are always explicit in Jimple. For the base cases, new taint are only produced at `AssignStmts`. We skip over `IdentityStmt` even though they also assign values such as the `this` reference or the parameters to locals. The Return Flow Function is responsible for mapping all parameters back into the callee and the return edge is after the `IdentityStmts`.

Now, lets consider the current statement is an `AssignStmt`. It consists of a variable, either a reference or a local, on the left side and an expression on the right side. The assignment has the structure  $x.f^n \leftarrow y.g^m$  with  $n, m \in \{0, 1\}$  modelling a possible field reference. As taints may have an access path of an arbitrary length, we denote the rest as  $h^k$ .<sup>1</sup> Jimple

---

<sup>1</sup> $h^k$  is a  $k$ -length access path, not  $k$ -times the same field reference

---

also ensures only one field dereference per statement which Arzt chose not to represent in the semi-formal definitions.

In the first case, we look at exact matches. Either we have an assignment with a local ( $n = 0$ ) or a field dereference ( $n = 1$ ). For both, the base variable needs to match. For the latter, also the first field of the access path has to match the field dereference. The first field dereference is removed from the taint and the remaining access path is copied to the newly created taint. The incoming taint is killed because it received the worthwhile information at this statement.

**Rule 1:** An incoming taint  $T = x.f^n.h^k$  with  $k \geq 0$  produces the outflowing taint set  $\{y.g^m.h^k\}$ .

Next, we need a rule for the wildcard. Recall section 2.3, symbolic access paths can also be  $k$ -limited to speedup the analysis and they are  $k = 5$ -limited in FLOWDROID by default. Thus, we might encounter a taint with no field dereferences and a wildcard appended. In this case, just the base needs to match. But this time, the left side is kept alive because we can not reason which field is tainted due to the cut-off approximation.

**Rule 2:** An incoming taint  $T = x.*$  with  $k \geq 0$  produces the outflowing taint set  $\{y.g^m.*, T\}$ .

Lastly, we could also observe a taint on the right side. In this case, we simply apply the identity, so propagate the taint over the statement. This case also gets important later when we consider aliasing in section 4.2.

**Rule 3:** An incoming taint  $T = y.g^m.h^k$  with  $k \geq 0$  produces the outflowing taint set  $\{T\}$ .

Whenever a taint neither matches on the left nor on the right side, we also perform the identity as the statement does not touch the contents of the tainted variable.

Rule 1 and Rule 3 work likewise with a wildcard appended. Similiary to the rest access path  $h^k$ , it is copied to the left side in Rule 1.

### 3.1.2. Call Flow

The Call Flow Function and subsequently also Return and Call TO Return Flow Function is applied whenever a statements contains an InvokeExpr. For call statements without assignments, we have statements of the structure  $o.m(a_0, \dots, a_n)$  with  $n \in \mathbb{N}$ .  $a_i$  denotes the  $i$ -th argument,  $p_i$  the  $i$ -th parameter and  $c$  the class instance of the callee's base object.



---

<pre> 1 void main() { 2     String str = "42"; 3     callee(str); 4     sink(str); 5 } 6 7 void callee(String str) { 8     str = source(); 9 } </pre>	<pre> 1 void main() { 2     SomeObject o = new         SomeObject(); 3     callee(o); 4     sink(o.str); 5 } 6 7 void callee(SomeObject o) { 8     o.str = source(); 9 } </pre>
(a) Taint Without Fields	(b) Taint With Fields

Figure 3.1.: Call Flow Example

When we encounter a tainted argument in the caller, the taint needs to go through the callee. Java uses pass-by-value, so the arguments are copied into the callee. Thus, for reference types, the pointer to the object is copied into the callee. If only the base reference is tainted but nothing more ( $k = 0$  and no wildcard), the callee can only access and overwrite the reference saved in the parameter on the stack but not change the data of the object on the heap. Due to the backwards direction, we know the update to such a taint can not be inside the callee. This gets more clear when we get specific to Java's types. Primitives do not have fields and strings are immutable<sup>2</sup>. Consider the example in Figure 3.1. On the left, we use the built-in String type. In line 2, `str` is copied into callee. At this point, `main` holds the `str` with the value "42" but pointing to another memory location than `str` in callee<sup>3</sup>. Thus, callee can write to `str` without changing the contents of `str` in `main`. On the right, the callee can update the field on the heap and therefore, the taint needs to be propagated into the callee to find the leak. Conclusively,  $k$  needs to be greater than 0.

**Rule 1:** An incoming taint  $T = a_i.h^k$  with  $k > 0 \wedge 0 \leq i \leq n$  produces the outflowing taint set  $\{p_i.h^k\}$ .

**Rule 2:** An incoming taint  $T = a_i.*$  with  $0 \leq i \leq n$  produces the outflowing taint set  $\{p_i.*\}$ .

---

<sup>2</sup>The immutable conditions results in transparent fields, e.g. we can treat strings as if they were primitives.

<sup>3</sup>The JVM might only set a copy-on-write flag on `str` in callee and point it to the same location as `str` in `main` to save memory. At least right before the update happens, it is guaranteed that the variable points to a different location.

---

For a non-static callee, it can also access instance fields of the base object. When we observe a tainted base object, the taint also needs to flow through the callee. The tainted object transforms into a *this* reference.

**Rule 3:** An incoming taint  $T = o.h^k$  with  $k > 0$  produces the outflowing taint set  $\{this_c.h^k\}$ .

**Rule 4:** An incoming taint  $T = o.*$  produces the outflowing taint set  $\{this_c.*\}$ .

Static fields are a special case. Their scope is the whole program and thus, are always visible in the callee. The taint is untouched as the reference is the same everywhere.

**Rule 5:** An incoming taint  $T = S.f.h^k$  with  $k \geq 0$  produces the outflowing taint set  $\{T\}$ .

In Jimple, an `AssignStmt` can also consist of a `Local` on the left side and an `InvokeExpr` on the right side. The structure of the statement is in this case  $x \leftarrow o.m(a_0, \dots, a_n)$ .  $r_i$  denotes a return value.  $m$  is the number of return statements in the callee. If we observe such a statement and the left side is tainted as well, we need to map the returned value into the callee. Now, methods can have multiple return statements and as we traverse the reversed interprocedural control-flow graph, there are multiple edges to go with. We can not reason which return statement is the right one, so we need to taint the operand of every return statement in the callee.

**Rule 6:** An incoming taint  $T = x.h^k$  with  $k \geq 0$  produces the outflowing taint set  $\{r_i.h^k \mid 0 \leq i < m\}$ .

Unlike at normal flows, we kill all taints not matching any of the rules. In the case of a taint being out of the callee's scope, the Call To Return flow function propagates the taint over the statement.

### 3.1.3. Return Flow

Reaching the end of a method, the taints need to be mapped back into the caller. The statement we consider is of the structure  $o.m(a_0, \dots, a_n)$  with  $n \in \mathbb{N}$ . Again,  $a_i$  denotes the  $i$ -th argument,  $p_i$  the  $i$ -th parameter and  $c$  the class instance.

The first rule matches rule 1 and 2 of call flow: map all parameters back into the caller. In contrast to call flow, this time we also map primitives and strings back into the caller because the contents of the variable were copied from the caller's argument into the callee.

---

**Rule 1:** An incoming taint  $T = p_i.h^k$  with  $k \geq 0 \wedge 0 \leq i \leq n$  produces the outflowing taint set  $\{a_i.h^k\}$ .

The *this* reference also needs to be mapped back into the caller.

**Rule 2:** An incoming taint  $t = this_c.h^k$  with  $k \geq 0$  produces the outflowing taint set  $T = \{o.h^k\}$ .

Tainted static fields are also mapped back untouched. This is the same as rule 5 of the Call Flow.

**Rule 3:** An incoming taint  $t = S.h^k$  with  $k \geq 0$  produces the outflowing taint set  $T = \{t\}$ .

Again, taints not matching any rule are killed. For example, this is the case when a method is left and the local variable is not in the scope of the caller.

### 3.1.4. CallToReturn Flow

The statement structure is  $o.m(a_0, \dots, a_n)$  with  $n \in \mathbb{N}$ .  $a_i$  denotes the  $i$ -th argument.

A taint is not in scope of the callee if it is not static and neither matches an argument nor the base object the method is called on. Such a taint is not matched inside Call Flow and needs to be propagated over the call statement.

**Rule 1:** An incoming taint  $T = x.h^k$  with  $k \geq 0 \wedge (\forall a \in Arguments : a \neq x) \wedge x \neq o \wedge x \notin Static$  produces the outflowing taint set  $T = \{t\}$ .

Consider again the left side of Figure 3.1. In line 3, the taint is in the kill set of Call Flow. But as we want to preserve the taint after the call, we need to propagate the taint over the call statement in such cases.

**Rule 2:** An incoming taint  $t = a_i$  with  $0 \leq i \leq n$  produces the outflowing taint set  $T = \{t\}$ .

Like in Call and Return Flow, here we also kill taints which do not match any of these rules.

---

## 3.2. Complexity of Data Flow Analysis

---

IFDS has a time-complexity of  $O(E \cdot D^3)$ . The edges in the control-flow graph are set by the to-be-analyzed app. The domain depends on the tainted variables observed by the IFDS analysis. Arzt et al evaluation of FLOWDROID shows no correlation between a-priori known parameters and the runtime of the analysis [2]. So, we are left with the number of taint propagations as the only parameter correlated to the runtime but they are only known afterwards.

The number of taint propagations depends on two factors: the lifetime of taints and the number of taints. Both factors highly depend on the search direction which we will explain in the following paragraphs.

```
1 String returnParam(int i, String s1, String s2, String s3) {
2     if (i == 1)
3         return s1;
4     else if (i == 2)
5         return s2;
6     else if (i == 3)
7         return s3;
8     else
9         return "default";
10 }
```

Figure 3.2.: Branching Factor Example

First, we take a look at the branching factor. The branching factor describes the number outgoing edges from a node. A smaller branching factor is favorable. Think of binary operator expression such as `int c = a + b;`, backwards we can not argue which operand is responsible for the tainted output and thus proceed with both operands tainted. The same restriction is present in rule 4 of Call Flow which describes how the returned value is mapped back into the callee. This time the branching factor can be even larger. As an example, in Figure 3.2 is a method which conditionally returns one of its parameters and is part of the leak path. Lets assume the returned value of a call to `returnParam()` is tainted. Backwards, every returned operand is tainted and later on mapped according to Return Flow rule 2 back into the caller. Thus, the IFDS algorithm ends up with a summary  $retVal \rightarrow \{s1, s2, s3\}$ . Forwards, a tainted parameter is mapped into the callee and later

on returned to the caller resulting in a summary  $sX \rightarrow \{retVal, sX\}$  with  $X \in \{1, 2, 3\}$ . Such a case favors forwards analysis.

In contrast, a strict right-to-left flow favors backwards analysis as backwards taints are killed more often due to a stronger overwrite rule meaning a shorter lifetime per taint. In Figure 3.3 is such a right-to-left flow displayed. Forwards, the right hand side is always kept alive because it still holds the tainted value below the statement and could be leaked. Starting at the sinks this fact can be confidently ignored because the fact is known to be not on a leak path.

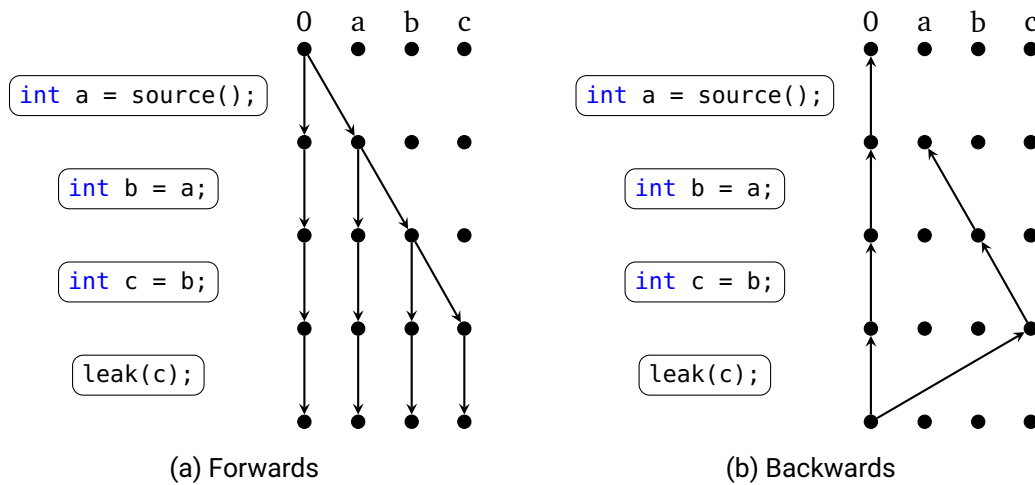


Figure 3.3.: Right-to-Left Order

Static field taints are a special case and an issue in both directions. Their scope is global, so a static field can be accessed anywhere in the code. Thus such a taint needs to be propagated into every method and can only be killed on overwrite.

In special cases of application, the direction of the dataflow analysis might be more important. We bring in classical web-oriented applications as an example. Typical vulnerabilities such as cross-site-scripting and SQL injection occur because untrusted user input was not validated. Normally the content is escaped right before output or sanitized before execution of the command, respectively. So the sanitization method is expected to be close to the sink. Being able to kill taints earlier is advantageous and therefore backwards analysis should be the better choice.



---

We discussed cases where one direction is better than the other but it still highly depends on the analyzed app.

- Number of taints
  - There seems to be no correlation between source count and analysis time
  - Probably also holds for sinks?
  - There might be indicator for a single app whether it is better to start at sources or sinks

---

## 4. Implementation

---

In this chapter, we describe the details of our backwards-directed implementation and how we integrated it into FLOWDROID.

---

### 4.1. Integration

---

FLOWDROID is built to be extensible from the ground up. We wanted to reuse as much components of FLOWDROID as possible.

First, we needed a backwards interprocedural control-flow graph. FLOWDROID already contained one for the on-demand aliasing which only missed the `notifyMethodChanged()` method. Next, we need to introduce unconditional taints at sinks and check for the matching access paths at sources. The methods for retrieving sources and sinks from a Source Sink Manager have different signatures because in the forwards case, only at sinks access paths have to match. We added the interface `IReversibleSourceSinkManager` extending the `ISourceSinkManager`. It enforces two additional methods:

- `SourceInfo getInverseSinkInfo(Stmt sCallSite, InfoflowManager manager)`
- `SinkInfo getInverseSourceInfo(Stmt sCallSite, InfoflowManager manager, AccessPath ap)`

`getInverseSinkInfo` returns the necessary information for introducing unconditional taints at sinks while `getInverseSourceInfo` also matches the access paths at sources. The source sink managers `DefaultSourceSinkManager` and `AccessPathBasedSourceSinkManager` now implement the `IReversibleSourceSinkManager` interface. Reversible source sink manager currently do not support the one-source-at-a-time mode.

For the core, we created two new classes implementing `IInfoflowProblem`. `BackwardsInfoflowProblem` implements the flow functions described in section 3.1. Additional language features are

---

sourced out into rules which are informally described in section 4.3. `BackwardsAliasProblem` is the on-demand forward alias analysis. We describe the on-demand aliasing in greater detail in section 4.2.

To hide the fact that we internally searched backwards, we also created a `BackwardsInfoflowResults` extending `InfoflowResults`. The implementation is quite simple. It overwrites only the `addResult` implementations to swap the source and sink. If full path reconstruction is enabled, we also reverse the path.

---

## 4.2. Flow-Sensitive Alias Analysis

---

FLOWDROID offers multiple aliasing strategies. In this work, we focus on the flow-sensitive alias analysis which is implemented as another IFDS problem called `BackwardsAliasProblem`. Basically, this is a forwards IFDS search with flow functions using aliasing rules.

```
1 void aliasRule1() {  
2     A a = b;  
3     b.str = source();  
4     sink(a.str);  
5 }
```

(a) Example for alias analysis initiated by rule 1

```
1 void aliasRule3() {  
2     A a = b;  
3     a.str = source();  
4     sink(b.str);  
5 }
```

(b) Example for alias analysis initiated by rule 3

Figure 4.1.: Normal flow Aliasing examples

**Handover** Whenever we visit a statement and notice a taint could have an alias, the taint is handed over to the alias analysis. Normal flow rule 3 is such a case. The taint is on the right side and we notice that the left side also refers to the same value in memory due to being stored in the heap. The left side gets tainted and propagated forwards to find out if we missed a write to the alias. In normal flow rule 1 and 2, we also turn around. Figure 4.1 shows two cases where the turnaround is necessary. In 4.1a, at line 2 `b` is assigned to `a`. Above this statement, the contents of `a` do not matter anymore, thus `a` is killed. On the other side, `b` is now tainted but as `a` and `b` aliases, we missed all updates to `b` below the statement. Concluding, at line 2 the `a.str` taint is killed, a new taint `b.str` is created and alias analysis for `b.str` is triggered. In 4.1b, we also observe that all writes to



---

the alias were neglected but this time the incoming taint is propagated over the statement and for `a.str` the alias analysis is triggered.

A taint is handed back to the infoflow search if the alias analysis encounters an update to the taint, e.g. the taint is on the left side of the assignment. Again consider Figure 4.1a. At line 3, `b.str` is tainted and on the left side of the assignment. At this point, `b.str` is handed over to the infoflow analysis with the statement on line 3. The backwards infoflow search then follows the missed path to detect possible leaks, in this case it right away reports a leak.

**Turn Unit** We added another field to the `Abstraction` class called `turnUnit`. This is the equivalent to the `activationUnit` in forwards analysis. The `turnUnit` references the last statement for which the taint is relevant in the infoflow search. At start, it is the sink where the taint was introduced. Later on, it is set whenever we visit an assignment with a primitive or string on the left side<sup>1</sup>. Consider the code in Figure 4.2. At line 5, the taint is introduced, line 3 taints `b.str` and sets the `turnUnit` to this statement. In line 2, `a` is found to be an alias of `b` and causes a handover to the alias problem. The `turnUnit` now stops the alias search at line 3 and prevents a false positive.

```
1 void turnStmtNeeded() {
2     A a = b;
3     String str = b.str;
4     a.str = source();
5     sink(str);
6 }
```

Figure 4.2.: Aliasing example with turn unit

---

## 4.3. Rules

---

Flow functions can get quite large, complicated to understand and hard to maintain [12]. To counteract this, `FLOWDROID` outsources certain features into rules. These rules also implement the four flow functions and are applied in the corresponding flow function. In this section, we describe our rule implementation and informally state the rule behavior.

---

<sup>1</sup>The `String` object is a special case in Java. It is saved on the heap but immutable and therefore can not alias.

---

### 4.3.1. Source & Sink Propagation Rule

In backwards analysis, sources act like sinks and vice versa. Thus, the Source Propagation Rule records taints flowing into sources and the Sink Propagation Rule unconditionally introduces taints at sinks.

Notably, the `DefaultSourceSinkManager` assumes the return value to be tainted and only if the return value is ignored or the method has no return value the base object is assumed to be tainted unless specified otherwise while at sinks base object and parameters are leaked [2]. Thus, starting at sinks results in more taints per start statement than in forwards analysis. As written in section 3.2, Arzt's evaluation has shown that the initial source count does not correlate with the runtime which implies that this should be insignificant.

### 4.3.2. Backwards Array Propagation Rule

The Array Propagation Rule handles `ArrayNewExpr`, `LengthExpr` and `ArrayRef` on the right hand side. Further, we describe the three cases of this rule.

- **Array Rule 1:** If the length of the left side is tainted and the right side is an `ArrayNewExpr`, the outcoming taint is the size local of the `ArrayNewExpr`.
- **Array Rule 2:** If the left side is tainted and the right side is a `LengthExpr`, the outcoming taint is the operand of the `LengthExpr` with only its length tainted.
- **Array Rule 3:** If the left side is tainted and the right side is an `ArrayRef`, the outcoming taint is the array base with only its content tainted.

By default, the whole array is tainted and indices are not tracked. Following, all three rules kill the incoming taint unless the left side is an `ArrayRef`.

### 4.3.3. Backwards Exception Propagation Rule

The Backwards Exception Propagation Rule handling is similar to forward.

We first observe a statement catching the exception such as `$someVar := @caughtexception`. If `someVar` is tainted, we derive a new taint with the throw flag set. In the next propagation of this taint, we either find the `ThrowStmt` in the normal flow or in the call flow. Similar

---

to call flow function, if the throw is inside another method, we need to iterate over all possible ThrowStmts in the callee.

#### 4.3.4. Backwards Wrapper Propagation Rule

FLOWDROID already provided a `IReversibleTaintWrapper` interface. Implementing taint wrappers support `getInverseTaints()` which takes the outgoing taint as an input and computes the incoming taints.

This rule is similar to its forward equivalent but enforces a reversible taint wrapper. Consequently, tainted return values are also passed into the taint wrapper.

Inverse taint wrappers do have one limitation. Often tainted parameters result in a tainted base object. `EasyTaintWrapper` uses this pattern to provide a fast and simple taint wrapper [2]. But backwards, only a tainted base object is observed. Similar to binary operators in assignments, we can not do anything but to taint every parameter.

#### 4.3.5. Backwards Strong Update Rule

Strong updates are assignments where the content of a variable gets overwritten. In our normal flow rules, this is modelled in rule 1. When a statement is observed with its left side tainted, we know it got its tainted content at this statement. Thus we kill the taint because the content above this statement is of no interest and taint the right hand side. So we performed a strong update on the left side.

But with aliasing this gets quite more complicated. Now, we can observe a taint not matching the left side and propagate it over the statement according to the default rule of normal flow but the taint is an alias of the left side and should have been killed. Linking aliasing taints to support such strong updates would lose the distributiveness property of the flow functions.

In this case, FLOWDROID falls back to Soot's must-aliasing analysis. However, the built-in must-aliasing is only intraprocedural. Thus, the strong update rule can not detect strong updates split over methods.

- **Strong Update Rule:** If the incoming taint must-aliases the left side then apply the normal flow rules just as if the left side was tainted.

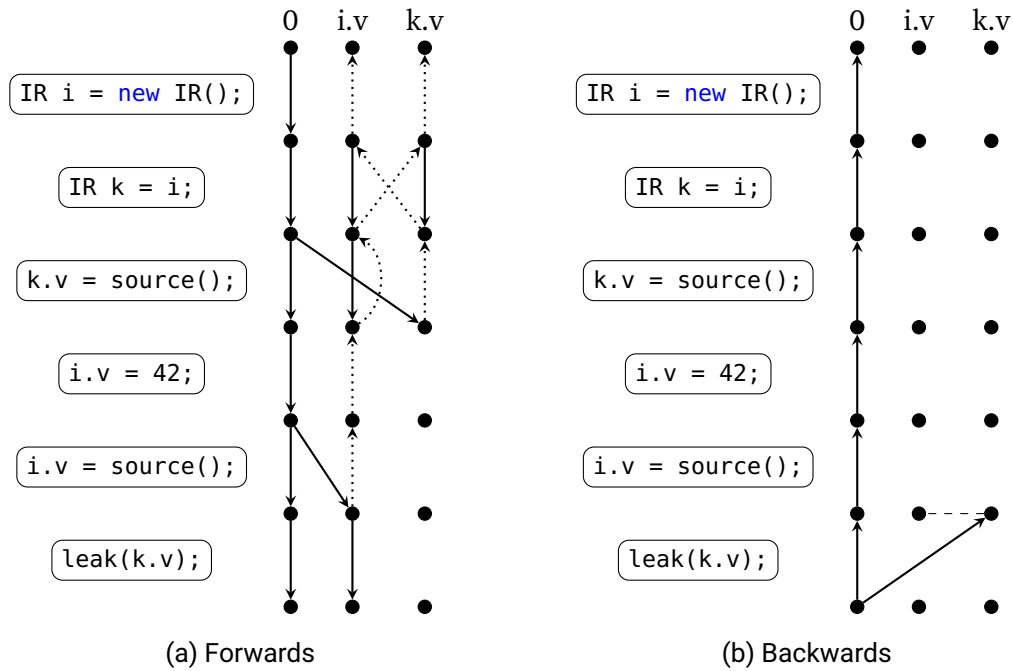


Figure 4.3.: Strong Update Example

Searching backwards also allows for easier strong updates. Backwards, the first observed strong update is the last before the sink.

Consider the example in Figure 4.3. IR is a wrapper around an integer to provide integers by reference. The main idea behind this example is a strong update on an alias of a tainted variable and later retain the alias. In 4.3a, taints are created at both source statements. The critical point is `i.v = 42;` where the forwards strong update rule kills `k.v` because it must alias `i.v` and the taint `i.v`, later found through aliasing, is killed according to normal flow. Both kills are appropriate because it is unknown whether there will be another strong update. This is not the case in backwards analysis. Here the last write before the sink to the leaked variable or one of its aliases is found first. So after one propagation the taint already reaches a source statement.

---

### 4.3.6. Backwards Clinit Rule

<clinit> is a special method in the JVM and stands for class loader init. The function is generated by the compiler and can not be called explicitly. Examples of statements which get compiled into clinit can be seen in Figure 4.4. The invocation is implicit at the initialization phase of the class and is executed at most once for each class<sup>2</sup>. This behavior is modelled as an overapproximation in FLOWDROID’s default call graph algorithm SPARK. SPARK adds an edge to <clinit> at each statement containing a StaticFieldRef, StaticInvokeExpr or NewExpr<sup>3</sup>.

<pre>1 class ClinitClass1 { 2     public static String str =         source(); 3 }</pre>	<pre>1 class ClinitClass2 { 2     static { 3         ClinitClass2.sink(); 4     } 5 }</pre>
(a) static variable initialization	(b) static block

Figure 4.4.: Examples of statements being in <clinit>

The need for this rule is rooted in the IFDS solver of FLOWDROID. The solver decides whether to use normal flow or call flow by calling `isCallStmt(Unit u)` on the interprocedural control-flow graph generated by Soot. Internally, this method calls `containsInvokeExpr()` on the Unit object. `containsInvokeExpr()` for `AssignStmt` only returns true if the right hand side is an instance of `InvokeExpr`. Consequently, the calls to <clinit> from `AssignStmts` with `NewExpr` or `StaticFieldRef` on the right side are missed.

The Backwards Clinit Rule manually injects an edge to the <clinit> method in the infowflow solver when appropriate during the analysis. Also, it lessens the overapproximation of SPARK by carefully choosing whether to inject the edge. The rule works as follows:

- **Clinit Rule 1:** If the tainted static variable is a field of the methods class: Do not inject because we will at least encounter a `NewExpr` of the same class further in the call graph.

---

<sup>2</sup><https://docs.oracle.com/javase/specs/jvms/se8/html/jvms-2.html#jvms-2.9>

<sup>3</sup><https://github.com/soot-oss/soot/blob/59931576784b910a7d38f81910b7313aa2feafea/src/main/java/soot/jimple/toolkits/callgraph/OnFlyCallGraphBuilder.java#L969>

- 
- **Clinit Rule 2:** Else if the tainted static variable matches the `StaticFieldRef` on the right hand side: Inject the edge because we can not be sure whether we see another edge to `<clinit>`.
  - **Clinit Rule 3:** Else if the class of the tainted static variable matches the class of the `NewExpr`: Inject the edge because we can not be sure whether we see another edge to `<clinit>`.

This is still an overapproximation of course. A precise solution would require bookkeeping of the first occurrence in the code of every class.

In forwards analysis, the issue is not as severe. As taints are introduced at sources, if the source statement is a static initialization as shown in Figure 4.4a, the propagation starts inside the `<clinit>` method. The solver has a `followReturnsPastSeeds` option which propagates return flows for unbalanced problems, for example when the taint was introduced inside a method and therefore there was no incoming flow. This allows the forwards analysis to detect leaks originated from static variable initializations but misses leaks inside static blocks as shown in Figure 4.4b.

#### 4.3.7. Other Rules

Skip System Class Rule and Stop After First K Flows Rule are not direction-dependent. Both are shared with the forwards search and therefore use the existing implementation in `FLOWDROID`.

---

### 4.4. Other Components

---

#### 4.4.1. Code Optimizer: AddNOPStmts

Before starting the analysis, `FLOWDROID` applies code optimization to the interprocedural call graph. By default, dead code elimination and within constant value propagation is performed. Those are also applied before backwards analysis but we needed another code optimizer to handle an edge case in backwards analysis.

First, take a look at `StaticTestCode#static2Test` in Figure 4.5. The method and entry point `static2Test` is static and does not have any parameters. Same is true for the source method `TelephonyManager#getDeviceId`. Due to the first condition, `static2Test` has

---

no identity statements and because of the second condition there are also no assign statements before the source statement in Jimple. Therefore the source statement is the first statement in the graph. Next, a detail of FLOWDROID's IFDS solver is important. The Return and CallToReturn flow function is only applied if a return site is available [2]. When searching backwards, the source statement is the last statement and thus has no return sites. Now recall subsection 4.3.1, taints flowing into sources are registered in the CallToReturn flow function. Altogether, leaks can not be found if the source statement is the first statement.

Moving the detection of incoming taints flows into sources from the CallToReturn to the Call flow function was not an option because by default source methods are not visited. Our solution is to just add a NOP statement in such cases. This saves us from introducing new edge cases inside the flow functions which are already complex enough. Due to the entry points being known beforehand, the overhead is negligible.

#### 4.4.2. Taint Wrappers

FLOWDROID already provides an interface `IReversibleTaintWrapper` for taint wrappers which also provide inversed summaries. The Summary Taint Wrapper for STUBDROID already implements this. For the Easy Taint Wrapper we contributed the inverse implementation.

#### 4.4.3. Native Call Handler

The native call handler of FLOWDROID handles two methods:

- `System#arraycopy`: If the first parameter is tainted, taint the second parameter.
- `reflect.Array#newArray`: If the length is tainted, propagate it over the statement.

We adapted the existing implementation and only reversed the logic of `System#arraycopy` to reflect the analysis direction.

---

```
1  public static void static2Test() {
2      String tainted = TelephonyManager.getDeviceId();
3      ClassWithStatic static1 = new ClassWithStatic();
4      static1.setTitle(tainted);
5      ClassWithStatic static2 = new ClassWithStatic();
6      String alsoTainted = static2.getTitle();
7
8      ConnectionManager cm = new ConnectionManager();
9      cm.publish(alsoTainted);
10 }
```

(a) Java

```
1  public static void static2Test() {
2      tainted = staticinvoke
          <soot.jimple.infoflow.test.android.TelephonyManager:
            java.lang.String getDeviceId()>(); // Line 2 in (a)
3
4      // [...]
5
6      virtualinvoke cm.<soot.jimple.infoflow.test.android.ConnectionManager:
          void publish(java.lang.String)>(alsoTainted); // Line 9 in (a)
7
8      return;
9  }
```

(b) Jimple

Figure 4.5.: static2Test Code



---

## 5. Validation

---

---

### 5.1. Unit Tests

---

FLOWDROID already contains 519 unit tests for the core component. We also validated the backwards analysis with these tests with good results. Because we focused on the context-sensitive analysis, not all analysis-related tests were applicable. In the following, we briefly explain why tests were left out or did not return the same results.

**EasyTaintWrapperTests** `equalsTest` and `hashCodeTest` are expected to return one leak but the backwards analysis does report no leaks. This difference is related to the `EasyTaintWrapper` implementation. The implementation marks `equals()` and `hashCode()` as exclusive. This means we can skip this method because we already have a rule for it. The check for exclusiveness is part of the `Call` and `CallToReturn` flow function. In both tests, the source is inside the `equals()` or `hashCode()` method. The IFDS solver behaves as already observed in subsection 4.3.6 and when searching forwards it creates a return edge returning from the method while going backwards we do not propagate into the method because it is exclusive. We marked those two tests forwards-specific and created two equivalent backwards-specific tests with sinks inside the `equals()` or `hashCode()` method with one expected leak.

**SourceSinkTests** These tests ensure the source sink manager can be swapped out. This is not relevant for the correctness of the backwards analysis and therefore are ignored.

**HeapTestPtsAliasing** We focused in this work on flow-sensitive aliasing. Points-To-Aliasing is left for future work.

**ImplicitFlowTests** Not implemented.

**SetTests** `containsTest` needs implicit flows to find the leak.

implement  
later if  
enough  
time

---

## 5.2. DroidBench

---

DROIDBENCH is a test suite to evaluate data flow analysis tools targeting the Android ecosystem. It originated from the initial work on FLOWDROID to assess it in comparison to other tools [3]. The latest development version 3 includes 190 test cases<sup>1</sup>. We used the newest commit on develop at the time of writing<sup>2</sup> for validating our implementation. We aim to achieve similar results as FLOWDROID's existing forwards implementation.

### 5.2.1. Configuration

For the validation we ran FLOWDROID with the default configuration of the Android module using the Easy Taint Wrapper as the taint wrapper. The configuration is summarized in Table 5.1.

Option	Value
Array Size Tainting	disabled
Inspect Sources & Sinks	disabled
Static Field Tracking	enabled
Ignore Flows in System Packages	enabled
Exclude Soot Library Classes	enabled
Timeout	-
Taint Wrapper	EasyTaintWrapper

Table 5.1.: Real World Apps Configuration

We only used a subset of DROIDBENCH's tests to validate our results. Dynamic Code Loading, Self Modification, Unreachable Code and Native Code are all not supported by FLOWDROID. The first three are all call-graph related and the latter is not supported because FLOWDROID has no Android native call handler for now. Also Inter Component Communication (ICC), Reflection Inter Component Communication and Inter App Communication were left out because the ICC module is - at the time of this work - not maintained anymore. All of the tests stated above are flow function independent. If FLOWDROID gets support for those features in the future they should also work in backward analysis.

---

<sup>1</sup><https://github.com/secure-software-engineering/DroidBench/>

<sup>2</sup>6th March 2021, Commit ddbd50c

## 5.2.2. Results

The full overview of the results is in Table 5.2.  $\odot$  denotes true positive,  $\star$  false positive and  $\bigcirc$  false negative. If a row is empty, the test expects no leaks and also none were found.

Our backwards-directed implementation yields nearly the same result as the existing forwards implementation with one missed leak more than the baseline. We achieve a F1 measure of 0.89 which is the same as the baseline.

App Name	Forwards	Backwards
Aliasing		
FlowSensitivity1		
Merge1	$\star$	$\star$
SimpleAliasing1	$\odot$	$\odot$
StrongUpdate1		
Arrays and Lists		
ArrayAccess1	$\star$	$\star$
ArrayAccess2	$\star$	$\star$
ArrayAccess3	$\odot$	$\odot$
ArrayAccess4		
ArrayAccess5		
ArrayCopy1	$\odot$	$\odot$
ArrayToString1	$\odot$	$\odot$
HashMapAccess1	$\star$	$\star$
ListAccess1	$\star$	$\star$
MultidimensionalArray1	$\odot$	$\odot$
Callbacks		
AnonymousClass1	$\odot$	$\odot$
Button1	$\odot$	$\odot$
Button2	$\odot \odot \odot \star$	$\odot \odot \odot \star$
Button3	$\odot \odot$	$\odot \odot$
Button4	$\odot$	$\odot$
Button5	$\odot$	$\odot$
LocationLeak1	$\odot \odot$	$\odot \odot$
LocationLeak2	$\odot \odot$	$\odot \odot$
LocationLeak3	$\odot$	$\odot$
MethodOverride1	$\odot$	$\odot$

App Name	Forwards	Backwards
MultiHandlers1		
Ordering1		
RegisterGlobal1	⊛	⊛
RegisterGlobal2	⊛	⊛
Unregister1	★	★
Emulator Detection		
Battery1	⊛	⊛
Bluetooth1	⊛	⊛
Build1	⊛	⊛
Contacts1	⊛	⊛
ContentProvider1	⊛⊛	⊛⊛
DeviceId1	⊛	⊛
File1	⊛	⊛
IMEI1	⊛⊛	○○
IP1	⊛	⊛
PI1	⊛	⊛
PlayStore1	⊛⊛	⊛⊛
PlayStore2	⊛	⊛
Sensors1	⊛	⊛
SubscriberId1	⊛	⊛
VoiceMail1	⊛	⊛
Field and Object Sensitivity		
FieldSensitivity1		
FieldSensitivity2		
FieldSensitivity3	⊛	⊛
FieldSensitivity4		
InheritedObjects1	⊛	⊛
ObjectSensitivity1		
ObjectSensitivity2		
Lifecycle		
ActivityEventSequence1	⊛	⊛
ActivityEventSequence2	○	○
ActivityEventSequence3	○	○
ActivityLifecycle1	⊛	⊛
ActivityLifecycle2	⊛	⊛
ActivityLifecycle3	⊛	⊛

App Name	Forwards	Backwards
ActivityLifecycle4	⊛	⊛
ActivitySavedState1	⊛	⊛
ApplicationLifecycle1	⊛	⊛
ApplicationLifecycle2	⊛	⊛
ApplicationLifecycle3	⊛	⊛
AsynchronousEventOrdering1	⊛	⊛
BroadcastReceiverLifecycle1	⊛	⊛
BroadcastReceiverLifecycle2	⊛ *	⊛ *
BroadcastReceiverLifecycle3	⊛	⊛
EventOrdering1	⊛	⊛
FragmentLifecycle1	⊛	⊛
FragmentLifecycle2	○	○
ServiceEventSequence1	○	○
ServiceEventSequence2	○	○
ServiceEventSequence3	○	○
ServiceLifecycle1	⊛	⊛
ServiceLifecycle2	⊛	⊛
SharedPreferencesChanged1	⊛ *	⊛ *
General Java		
Clone1	⊛	⊛
Exceptions1	⊛	⊛
Exceptions2	⊛	⊛
Exceptions3	*	*
Exceptions4	⊛	⊛
Exceptions5	⊛	⊛
Exceptions6	⊛	⊛
Exceptions7		
FactoryMethods1	⊛ ⊛	⊛ ⊛
Loop1	⊛	⊛
Loop2	⊛	⊛
Serialization1	○	○
SourceCodeSpecific1	⊛	⊛
StartProcessWithSecret1	⊛	⊛
StaticInitialization1	○	⊛
StaticInitialization2	⊛	⊛
StaticInitialization3	○	○

App Name	Forwards	Backwards
StringFormatter1	○	○
StringPatternMatching1	★	★
StringToCharArray1	★	★
StringToOutputStream1	★ ★	★ ★
UnreachableCode		
VirtualDispatch1	★ ★	★ ★
VirtualDispatch2	★ ★	★ ★
VirtualDispatch3	★	★
VirtualDispatch4		
Miscellaneous Android-Specific		
ApplicationModeling1	★	★
DirectLeak1	★	★
InactiveActivity		
Library2	★	★
LogNoLeak		
Obfuscation1	★	★
Parcel1	★	★
PrivateDataLeak1	★	★
PrivateDataLeak2	★	★
PrivateDataLeak3	★ ○	★ ○
PublicAPIField1	★	★
PublicAPIField2	★	★
View1	★	★
Reflection		
Reflection1	★	★
Reflection2	★	★
Reflection3	★	★
Reflection4	★	★
Reflection5	★	★
Reflection6	★	★
Reflection7	○	○
Reflection8	★	★
Reflection9	★	★
Threading		
AsyncTask1	★	★
Executor1	★	★

App Name	Forwards	Backwards
JavaThread1	⊛	⊛
JavaThread2	⊛	⊛
Looper1	⊛	⊛
TimerTask1	⊛	⊛
⊛	103	102
★	13	13
○	12	13
Precision	88.79%	88.7%
Recall	89.57%	88.7%
F1 measure	0.89	0.89

Table 5.2.: DroidBench Validation Results

### 5.2.3. Results Explanation

In greater detail, we miss both leaks in `EmulatorDetectionTests#IMEI1` whereas in `StaticTests#StaticInitialization1` we do not miss the leak. We provide explanation on why below.

#### Emulator Detection

**IMEI1** needs implicit flows to find both leaks. The source is only used inside the condition of an if statement. See Figure 5.1, we are unable to taint the variable `imei` without an implicit taint created in line 8.

could  
be fixed  
in the  
future

#### General Java

As all `StaticInitialization` tests depends on the modelling of the `<clinit>` behavior, we decided to explain all three even though only `StaticInitialization1` is different.

**StaticInitialization1** differs in forwards and backwards analysis. Backwards it reports one leak due to the explicit modelling of `<clinit>` edges instead of relying on SPARK. Recall subsection 4.3.6, leaks inside static blocks are missed in forward analysis. This test

---

```

1 String imei = telephonyManager.getDeviceId(); // source
2 String suffix = "0000000000000000"; // T={}
3 String prefix = "secret"; // T={}
4 String msg = prefix + suffix; // T={prefix, suffix}
5
6 int zeroPos = 0; // zeroPos dies here
7 while (zeroPos < imei.length()) {
8     if (imei.charAt(zeroPos) == '0') // implicit flow needed
9         zeroPos++;
10    else {
11        zeroPos = 0;
12        break;
13    }
14 }
15
16 String newImei = msg.substring(zeroPos, zeroPos + Math.min(prefix.length(),
    msg.length() - 1)); // T={msg, zeroPos}
17 Log.d("DROIDBENCH", newImei); // T={newImei}
18
19 SmsManager sm = SmsManager.getDefault();
20 sm.sendTextMessage("+49 123", null, newImei, null, null); // T={newImei}

```

Figure 5.1.: IMEI1 excerpt

case is quite similar to Figure 4.4b and therefore the leak is only reported in backwards analysis at the moment.

**StaticInitialization2** yields the same result but because of different reasons. The test assigns a tainted value to a static field in the static initializer. Again, recall subsection 4.3.6. Backwards, the `clinit` rule takes care of visiting the `<clinit>` edge while forwards the `followReturnsPastSeeds` option of the IFDS solver is responsible.

**StaticInitialization3**'s leak is missed despite the explicit modelling of `clinit`. The code is provided in Figure 5.2. `MainActivity` is using the singleton pattern and thus has a static field `v` referring to its instance. The source statement is inside the static block of the `Test` class using the singleton to access the instance field `s`. The taint is now introduced at the sink and refers to the field through the `this` instance. When we visit line 13, the



---


<clinit> edge is not taken due to the taint being an instance field. Line 12 kills the only taint and stops the analysis as there is no taint to propagate anymore. We never get to see the statement where the static field `v` aliases `this`. This is a limitation of the alias analysis where only encountered aliases are tracked. Tracking more aliases is too inefficient for the analysis to be applicable and also fully precise aliasing is still an unsolved problem [4].

```
1 public class MainActivity extends Activity {
2     public static MainActivity v;
3     public String s;
4
5     @Override
6     protected void onCreate(Bundle savedInstanceState) {
7         v = this;
8
9         super.onCreate(savedInstanceState);
10        setContentView(R.layout.activity_main);
11
12        s = ""; // T={}
13        Test t = new Test(); // T={this.s}
14        Log.i("DroidBench", s); // T={this.s}
15    }
16 }
17
18 class Test {
19     static {
20         TelephonyManager mgr = (TelephonyManager)
21             MainActivity.v.getSystemService(Activity.TELEPHONY_SERVICE);
22         MainActivity.v.s = mgr.getDeviceId(); // source
23     }
24 }
```

Figure 5.2.: StaticInitialization3 code

#### 5.2.4. Improvements From The Summary Taint Wrapper

We briefly explained the simple but not always precise rules of the Easy Taint Wrapper in



---

subsection 4.4.2. Using STUBDROID's more precise summaries yields even better results for both directions. The false positives in the test cases `BroadcastReceiverLifecycle2` and `SharedPreferencesChanged1` are gone and the leak in `Serialization1` is found.

---

## 6. Performance Evaluation

---

In the last chapter, we have shown that our implementation has the necessary soundness to be viable and yields the expected results. We now evaluate our implementation against the existing implementation in FLOWDROID.

---

### 6.1. DroidBench

---

We already introduced DROIDBENCH in section 5.2 to validate the soundness of our backward-directed implementation. In this section, we focus on the performance in comparison to the existing forward-directed implementation in FLOWDROID.

DROIDBENCH has the advantage that all apps are crafted explicitly for benchmarking taint analysis. So, most tests only contain a single-figure number of sources and sinks. Also, the number of sources and sinks are often equal or differ by one to test whether the tool can differentiate something. These simplify the comparison between both analysis directions as neither one has an initial disadvantage.

Most test cases are small enough to be analyzed in sub-two seconds on an average four-core desktop CPU from 2012. Our test environment is not isolated, so background tasks and the process scheduler can affect the runtime. The short runtime, together with the variance of the unisolated testing environment, render the runtime unusable as a comparison point. In contrast, edge propagations are deterministic<sup>1</sup> and correlate with the runtime. Thus, we only use the number of propagations to compare both implementations.

The configuration is the same as described in subsection 5.2.1.

---

<sup>1</sup>This is only true if there are enough resources. FLOWDROID tries to gracefully terminate when running low on memory. Also, timeouts result in a non-reproducible number of edge propagations.

### 6.1.1. Results

We compare all test cases where both implementations yield the same result. When rows only contain hyphens, either the result of the test case differed between the two analyses or the IFDS analysis did not start, e.g., because no sink is in the reachable code. #I denotes the number of edge propagations inside the infoflow analysis and #A the number of edge propagations inside the alias analysis. We calculated the absolute difference as  $Result_B - Result_F$ . The relative difference calculates as follows:  $\frac{TotalDifference}{|\#I_F + \#A_F|}$ . Hence, negative values signify the backward analysis performed better. The full results are in Table 6.1.

In general, both implementations have similar average edge propagation counts. There are not many test cases where both perform identically; instead, dependent on the specific test case, the relative difference is between  $-1$  and  $1$ . So, the expected behavior from section 3.2 occurred: it highly depends on the analyzed app. However, we did not expect cases that let the backward edge propagations explode up to a factor of 10000%, as seen in `LifecycleTest#BroadcastReceiverLifecycle3` and others. In contrast, the existing forward implementation only at most a relative difference of 100%.

App Name	Forwards		Backwards		Difference			
	#I	#A	#I	#A	#I	#A	Total	Relative
AliasingTest								
FlowSensitivity1	175	72	39	4	-136	-68	-204	-0.83
Merge1	94	44	61	9	-33	-35	-68	-0.49
SimpleAliasing1	35	13	20	3	-15	-10	-25	-0.52
StrongUpdate1	30	13	11	3	-19	-10	-29	-0.67
AndroidSpecificTest								
ApplicationModeling1	235	103	851	1208	616	1105	1721	5.09
DirectLeak1	3	0	4	0	1	0	1	0.33
InactiveActivity	—	—	—	—	—	—	—	—
Library2	5	0	6	0	1	0	1	0.2
LogNoLeak	—	—	—	—	—	—	—	—
Obfuscation1	4	0	4	0	0	0	0	0.0
Parcel1	144	15	66	68	-78	53	-25	-0.16
PrivateDataLeak1	410	110	599	835	189	725	914	1.76
PrivateDataLeak2	15	0	5	6	-10	6	-4	-0.27
PrivateDataLeak3	17	2	212	143	195	141	336	17.68
runPublicAPIField1	89	1	62	31	-27	30	3	0.03
runPublicAPIField2	5	0	11	1	6	1	7	1.4
runView1	71	50	69	0	-2	-50	-52	-0.43
ArrayAndListTest								

App Name	Forwards		Backwards		Difference			
	#I	#A	#I	#A	#I	#A	Total	Relative
ArrayAccess1	77	34	51	100	-26	66	40	0.36
ArrayAccess2	16	4	12	0	-4	-4	-8	-0.4
ArrayAccess3	77	34	51	100	-26	66	40	0.36
ArrayAccess4	164	84	42	21	-122	-63	-185	-0.75
ArrayAccess5	75	5	67	63	-8	58	50	0.62
ArrayCopy1	18	2	9	2	-9	0	-9	-0.45
ArrayToString1	10	1	6	1	-4	0	-4	-0.36
HashMapAccess1	22	5	15	1	-7	-4	-11	-0.41
ListAccess1	85	9	60	97	-25	88	63	0.67
MultidimensionalArray1	29	3	16	23	-13	20	7	0.22
CallbackTest								
AnonymousClass1	152	0	208	1	56	1	57	0.38
Button1	58	39	43	0	-15	-39	-54	-0.56
Button2	454	66	155	257	-299	191	-108	-0.21
Button3	355	89	109	408	-246	319	73	0.16
Button4	58	39	43	0	-15	-39	-54	-0.56
Button5	80	40	6	6	-74	-34	-108	-0.9
LocationLeak1	617	222	260	300	-357	78	-279	-0.33
LocationLeak2	212	121	152	2	-60	-119	-179	-0.54
LocationLeak3	259	73	104	117	-155	44	-111	-0.33
MethodOverride1	3	0	2	0	-1	0	-1	-0.33
MultiHandlers1	17	0	145	151	128	151	279	16.41
Ordering1	456	151	44	2	-412	-149	-561	-0.92
RegisterGlobal1	207	103	49	0	-158	-103	-261	-0.84
RegisterGlobal2	52	37	43	0	-9	-37	-46	-0.52
Unregister1	11	0	9	1	-2	1	-1	-0.09
EmulatorDetectionTest								
Battery1	7	0	43	15	36	15	51	7.29
Bluetooth1	4	0	4	0	0	0	0	0.0
Build1	4	0	4	0	0	0	0	0.0
Contacts1	52	0	210	19	158	19	177	3.4
ContentProvider1	13	0	8	0	-5	0	-5	-0.38
DeviceId1	15	0	6	0	-9	0	-9	-0.6
File1	4	0	4	0	0	0	0	0.0
IMEI1	129	0	140	34	11	34	45	0.35
IP1	4	0	29	1	25	1	26	6.5
PI1	6	0	4	0	-2	0	-2	-0.33
PlayStore1	158	0	8	0	-150	0	-150	-0.95
PlayStore2	4	0	4	0	0	0	0	0.0
Sensors1	5	0	4	0	-1	0	-1	-0.2
SubscriberId1	29	0	4	0	-25	0	-25	-0.86
VoiceMail1	4	0	4	0	0	0	0	0.0
FieldAndObjectSensitivityTest								
FieldSensitivity1	98	50	25	3	-73	-47	-120	-0.81

App Name	Forwards		Backwards		Difference			
	#I	#A	#I	#A	#I	#A	Total	Relative
FieldSensitivity2	35	15	19	0	-16	-15	-31	-0.62
FieldSensitivity3	38	15	16	0	-22	-15	-37	-0.7
FieldSensitivity4	14	6	8	0	-6	-6	-12	-0.6
InheritedObjects1	4	0	6	0	2	0	2	0.5
ObjectSensitivity1	19	7	14	1	-5	-6	-11	-0.42
ObjectSensitivity2	15	8	10	0	-5	-8	-13	-0.57
GeneralJavaTest								
Clone1	23	2	12	4	-11	2	-9	-0.36
Exceptions1	16	0	13	0	-3	0	-3	-0.19
Exceptions2	22	0	13	0	-9	0	-9	-0.41
Exceptions3	18	0	11	0	-7	0	-7	-0.39
Exceptions4	20	1	22	1	2	0	2	0.1
Exceptions5	13	1	16	1	3	0	3	0.21
Exceptions6	77	12	23	0	-54	-12	-66	-0.74
Exceptions7	71	12	6	0	-65	-12	-77	-0.93
FactoryMethods1	40	0	14	2	-26	2	-24	-0.6
Loop1	93	2	46	7	-47	5	-42	-0.44
Loop2	123	2	74	14	-49	12	-37	-0.3
Serialization1	50	4	22	29	-28	25	-3	-0.06
SourceCodeSpecific1	16	0	45	7	29	7	36	2.25
StartProcessWithSecret1	29	8	17	3	-12	-5	-17	-0.46
StaticInitialization1	-	-	9	0	-	-	-	-
StaticInitialization2	57	29	86	0	29	-29	0	0.0
StaticInitialization3	35	9	5	0	-30	-9	-39	-0.89
StringFormatter1	16	1	10	1	-6	0	-6	-0.35
StringPatternMatching1	23	1	8	6	-15	5	-10	-0.42
StringToCharArray1	91	4	42	6	-49	2	-47	-0.49
StringToOutputStream1	26	3	30	3	4	0	4	0.14
UnreachableCode	-	-	-	-	-	-	-	-
VirtualDispatch1	128	31	88	28	-40	-3	-43	-0.27
VirtualDispatch2	7	0	12	0	5	0	5	0.71
VirtualDispatch3	8	0	6	0	-2	0	-2	-0.25
VirtualDispatch4	-	-	-	-	-	-	-	-
LifecycleTest								
ActivityEventSequence1	58	35	73	0	15	-35	-20	-0.22
ActivityEventSequence2	32	24	77	0	45	-24	21	0.38
ActivityEventSequence3	233	116	156	1	-77	-115	-192	-0.55
ActivityLifecycle1	99	72	156	7	57	-65	-8	-0.05
ActivityLifecycle2	47	34	33	0	-14	-34	-48	-0.59
ActivityLifecycle3	65	31	28	0	-37	-31	-68	-0.71
ActivityLifecycle4	49	33	14	0	-35	-33	-68	-0.83
ActivitySavedState1	20	0	7	1	-13	1	-12	-0.6
ApplicationLifecycle1	37	10	82	0	45	-10	35	0.74
ApplicationLifecycle2	86	17	94	155	8	138	146	1.42

App Name	Forwards		Backwards		Difference			
	#I	#A	#I	#A	#I	#A	Total	Relative
ApplicationLifecycle3	32	12	21	0	-11	-12	-23	-0.52
AsynchronousEventOrdering1	51	31	16	0	-35	-31	-66	-0.8
BroadcastReceiverLifecycle1	4	0	4	0	0	0	0	0.0
BroadcastReceiverLifecycle2	109	44	248	114	139	70	209	1.37
BroadcastReceiverLifecycle3	3	0	195	110	192	110	302	100.67
EventOrdering1	61	29	30	0	-31	-29	-60	-0.67
FragmentLifecycle1	187	127	90	0	-97	-127	-224	-0.71
FragmentLifecycle2	—	—	—	—	—	—	—	—
ServiceEventSequence1	53	20	152	34	99	14	113	1.55
ServiceEventSequence2	64	21	389	220	325	199	524	6.16
ServiceEventSequence3	46	12	275	151	229	139	368	6.34
ServiceLifecycle1	119	44	42	0	-77	-44	-121	-0.74
ServiceLifecycle2	68	20	89	21	21	1	22	0.25
SharedPreferenceChanged1	13	0	20	1	7	1	8	0.62
ReflectionTest								
Reflection1	15	5	8	0	-7	-5	-12	-0.6
Reflection2	21	5	11	0	-10	-5	-15	-0.58
Reflection3	42	9	62	25	20	16	36	0.71
Reflection4	9	0	8	0	-1	0	-1	-0.11
Reflection5	16	1	11	0	-5	-1	-6	-0.35
Reflection6	7	0	134	51	127	51	178	25.43
Reflection7	15	5	15	11	0	6	6	0.3
Reflection8	35	7	14	0	-21	-7	-28	-0.67
Reflection9	42	7	21	0	-21	-7	-28	-0.57
ThreadingTest								
AsyncTask1	22	2	11	1	-11	-1	-12	-0.5
Executor1	34	7	17	0	-17	-7	-24	-0.59
JavaThread1	34	7	17	0	-17	-7	-24	-0.59
JavaThread2	62	10	31	8	-31	-2	-33	-0.46
Looper1	49	3	20	16	-29	13	-16	-0.31
TimerTask1	203	28	32	33	-171	5	-166	-0.72
∅ Propagations	70.74	21.42	61.94	41.54	-8.8	20.11	11.32	1.41

Table 6.1.: DroidBench Performance Evaluation Results

### 6.1.2. Result Explanation

We define tests with a relative difference greater than 10 as worth investigating. In the following, we explain why our implementation performed worse than expected.

---

**PrivateDataLeak3** This test contains two sinks and one source. The tainted data is written to a file, later read from the file and then leaked. FLOWDROID does not support tracking taints over files, so it only finds a leak from source to file write but misses the leak from file read to send SMS. Due to EasyTaintWrapper's simplicity, overtainting happens in the backward direction. When `FileInputStream fis = openFileInput("out.txt");` is called with `fis` tainted, EasyTaintWrapper also taints the base object - the `MainActivity` in this case. As the `MainActivity` has a enormous scope, the taint has a long lifetime and many other taints could derive from this taint. This taint explains the relative difference of 17.68. Using the more precise SummaryTaintWrapper, the edges reduce to (51, 16) and a relative difference of 2.53, which is more reasonable. It is still higher because of the second sink.

**MultiHandlers1** Two `LocationListeners` are registered in different activities. In both activities, an instance field is a parameter of a sink. So there are two possible paths where something could be leaked. The `LocationListener` does not call any source on the first path, while the second path has an empty setter method killing the taint. For the first path, the backward analysis has to propagate the taint into the `LocationListener` to notice that this is a dead-end while the forward's search does not even start there. For the second path, the backward analysis seems to suffer because it starts at an instance field taint with a larger scope than a local variable.

**BroadcastReceiverLifecycle3** The test contains five sinks but only one source. If we only consider the leak path, both implementations perform equally. The four other sinks are responsible for the overhead on edge propagations.

**Reflection6** The reflective call site has multiple callees in the interprocedural control-flow graph. Backward all of these callees are visited, of which only one contains a source statement. Forwards, the taint is introduced in the callee at the source and just one return site needs to be processed.

### 6.1.3. Using A More Precise Taint Wrapper

We noticed the overtainting in `PrivateDataLeak3` is caused by the Easy Taint Wrapper, thus we now compare how using the Summary Taint Wrapper lowers the edge propagations.



As we already described, PrivateDataLeak3 benefits from the more precise taint wrapper.

App Name	EasyTW		SummaryTW		Difference			
	#I	#A	#I	#A	#I	#A	Total	Relative
AliasingTest								
FlowSensitivity1	39	4	71	13	32	9	41	0.95
Merge1	61	9	109	91	48	82	130	1.86
SimpleAliasing1	20	3	20	3	0	0	0	0.0
StrongUpdate1	11	3	11	3	0	0	0	0.0
AndroidSpecificTest								
ApplicationModeling1	851	1208	427	792	-424	-416	-840	-0.41
DirectLeak1	4	0	4	0	0	0	0	0.0
InactiveActivity	-	-	-	-	-	-	-	-
Library2	6	0	6	0	0	0	0	0.0
LogNoLeak	-	-	-	-	-	-	-	-
Obfuscation1	4	0	4	0	0	0	0	0.0
Parcel1	66	68	87	76	21	8	29	0.22
PrivateDataLeak1	599	835	576	835	-23	0	-23	-0.02
PrivateDataLeak2	5	6	5	6	0	0	0	0.0
PrivateDataLeak3	212	143	41	16	-171	-127	-298	-0.84
runPublicAPIField1	62	31	55	16	-7	-15	-22	-0.24
runPublicAPIField2	11	1	14	1	3	0	3	0.25
runView1	69	0	69	0	0	0	0	0.0
ArrayAndListTest								
ArrayAccess1	51	100	51	100	0	0	0	0.0
ArrayAccess2	12	0	12	0	0	0	0	0.0
ArrayAccess3	51	100	51	100	0	0	0	0.0
ArrayAccess4	42	21	42	21	0	0	0	0.0
ArrayAccess5	67	63	67	63	0	0	0	0.0
ArrayCopy1	9	2	9	2	0	0	0	0.0
ArrayToString1	6	1	6	1	0	0	0	0.0
HashMapAccess1	15	1	15	1	0	0	0	0.0
ListAccess1	60	97	77	118	17	21	38	0.24
MultidimensionalArray1	16	23	16	23	0	0	0	0.0
CallbackTest								
AnonymousClass1	208	1	208	1	0	0	0	0.0
Button1	43	0	43	0	0	0	0	0.0
Button2	155	257	184	275	29	18	47	0.11
Button3	109	408	120	357	11	-51	-40	-0.08
Button4	43	0	43	0	0	0	0	0.0
Button5	6	6	7	7	1	1	2	0.17
LocationLeak1	260	300	286	316	26	16	42	0.07
LocationLeak2	152	2	152	2	0	0	0	0.0
LocationLeak3	104	117	107	117	3	0	3	0.01
MethodOverride1	2	0	2	0	0	0	0	0.0
MultiHandlers1	145	151	148	151	3	0	3	0.01



App Name	EasyTW		SummaryTW		Difference			
	#I	#A	#I	#A	#I	#A	Total	Relative
Ordering1	44	2	44	2	0	0	0	0.0
RegisterGlobal1	49	0	49	0	0	0	0	0.0
RegisterGlobal2	43	0	43	0	0	0	0	0.0
Unregister1	9	1	9	1	0	0	0	0.0
EmulatorDetectionTest								
Battery1	43	15	39	15	-4	0	-4	-0.07
Bluetooth1	4	0	4	0	0	0	0	0.0
Build1	4	0	4	0	0	0	0	0.0
Contacts1	210	19	167	4	-43	-15	-58	-0.25
ContentProvider1	8	0	8	0	0	0	0	0.0
DeviceId1	6	0	6	0	0	0	0	0.0
File1	4	0	4	0	0	0	0	0.0
IMEI1	140	34	111	74	-29	40	11	0.06
IP1	29	1	54	6	25	5	30	1.0
PI1	4	0	4	0	0	0	0	0.0
PlayStore1	8	0	8	0	0	0	0	0.0
PlayStore2	4	0	4	0	0	0	0	0.0
Sensors1	4	0	4	0	0	0	0	0.0
SubscriberId1	4	0	4	0	0	0	0	0.0
VoiceMail1	4	0	4	0	0	0	0	0.0
FieldAndObjectSensitivityTest								
FieldSensitivity1	25	3	25	3	0	0	0	0.0
FieldSensitivity2	19	0	19	0	0	0	0	0.0
FieldSensitivity3	16	0	16	0	0	0	0	0.0
FieldSensitivity4	8	0	8	0	0	0	0	0.0
InheritedObjects1	6	0	6	0	0	0	0	0.0
ObjectSensitivity1	14	1	14	1	0	0	0	0.0
ObjectSensitivity2	10	0	10	0	0	0	0	0.0
GeneralJavaTest								
Clone1	12	4	19	10	7	6	13	0.81
Exceptions1	13	0	13	0	0	0	0	0.0
Exceptions2	13	0	13	0	0	0	0	0.0
Exceptions3	11	0	11	0	0	0	0	0.0
Exceptions4	22	1	53	26	31	25	56	2.43
Exceptions5	16	1	29	6	13	5	18	1.06
Exceptions6	23	0	23	0	0	0	0	0.0
Exceptions7	6	0	6	0	0	0	0	0.0
FactoryMethods1	14	2	14	2	0	0	0	0.0
Loop1	46	7	46	7	0	0	0	0.0
Loop2	74	14	74	14	0	0	0	0.0
Serialization1	22	29	332	547	310	518	828	16.24
SourceCodeSpecific1	45	7	45	7	0	0	0	0.0
StartProcessWithSecret1	17	3	18	4	1	1	2	0.1
StaticInitialization1	9	0	9	0	0	0	0	0.0

App Name	EasyTW		SummaryTW		Difference			
	#I	#A	#I	#A	#I	#A	Total	Relative
StaticInitialization2	86	0	85	0	-1	0	-1	-0.01
StaticInitialization3	5	0	5	0	0	0	0	0.0
StringFormatter1	10	1	10	1	0	0	0	0.0
StringPatternMatching1	8	6	7	1	-1	-5	-6	-0.43
StringToCharArray1	42	6	42	6	0	0	0	0.0
StringToOutputStream1	30	3	26	3	-4	0	-4	-0.12
UnreachableCode	-	-	-	-	-	-	-	-
VirtualDispatch1	88	28	110	88	22	60	82	0.71
VirtualDispatch2	12	0	12	0	0	0	0	0.0
VirtualDispatch3	6	0	6	0	0	0	0	0.0
VirtualDispatch4	-	-	-	-	-	-	-	-
LifecycleTest								
ActivityEventSequence1	73	0	73	0	0	0	0	0.0
ActivityEventSequence2	77	0	77	0	0	0	0	0.0
ActivityEventSequence3	156	1	156	1	0	0	0	0.0
ActivityLifecycle1	156	7	156	7	0	0	0	0.0
ActivityLifecycle2	33	0	33	0	0	0	0	0.0
ActivityLifecycle3	28	0	28	0	0	0	0	0.0
ActivityLifecycle4	14	0	14	0	0	0	0	0.0
ActivitySavedState1	7	1	7	1	0	0	0	0.0
ApplicationLifecycle1	82	0	82	0	0	0	0	0.0
ApplicationLifecycle2	94	155	94	155	0	0	0	0.0
ApplicationLifecycle3	21	0	21	0	0	0	0	0.0
AsynchronousEventOrdering1	16	0	16	0	0	0	0	0.0
BroadcastReceiverLifecycle1	4	0	4	0	0	0	0	0.0
BroadcastReceiverLifecycle2	248	114	208	98	-40	-16	-56	-0.15
BroadcastReceiverLifecycle3	195	110	144	82	-51	-28	-79	-0.26
EventOrdering1	30	0	30	0	0	0	0	0.0
FragmentLifecycle1	90	0	90	0	0	0	0	0.0
FragmentLifecycle2	-	-	-	-	-	-	-	-
ServiceEventSequence1	152	34	122	38	-30	4	-26	-0.14
ServiceEventSequence2	389	220	315	176	-74	-44	-118	-0.19
ServiceEventSequence3	275	151	232	110	-43	-41	-84	-0.2
ServiceLifecycle1	42	0	42	0	0	0	0	0.0
ServiceLifecycle2	89	21	89	21	0	0	0	0.0
SharedPreferenceChanged1	20	1	8	0	-12	-1	-13	-0.62
ReflectionTest								
Reflection1	8	0	8	0	0	0	0	0.0
Reflection2	11	0	11	0	0	0	0	0.0
Reflection3	62	25	49	0	-13	-25	-38	-0.44
Reflection4	8	0	8	0	0	0	0	0.0
Reflection5	11	0	11	0	0	0	0	0.0
Reflection6	134	51	122	31	-12	-20	-32	-0.17
Reflection7	15	11	3	0	-12	-11	-23	-0.88

App Name	EasyTW		SummaryTW		Difference			
	#I	#A	#I	#A	#I	#A	Total	Relative
Reflection8	14	0	14	0	0	0	0	0.0
Reflection9	21	0	21	0	0	0	0	0.0
ThreadingTest								
AsyncTask1	11	1	11	1	0	0	0	0.0
Executor1	17	0	17	0	0	0	0	0.0
JavaThread1	17	0	17	0	0	0	0	0.0
JavaThread2	31	8	28	8	-3	0	-3	-0.08
Looper1	20	16	20	16	0	0	0	0.0
TimerTask1	32	33	45	37	13	4	17	0.26
∅ Propagations	61.52	41.2	58.44	41.27	-3.07	0.06	-3.01	0.17
∅ without Serialization1	61.84	41.3	56.22	37.15	-5.62	-4.15	-9.76	0.04

Table 6.2.: DROIDBENCH Evaluation with Summary Taint Wrapper

## 6.2. Real World Apps

### 6.2.1. Configuration

Our test machine is equipped with four Intel Xeon E5-4650 and 1 TB of RAM. We limited the JVM to 50 GB RAM and FLOWDROID on 16 threads per instance. We ran at most four instances in parallel to ensure a one-to-one mapping between CPU threads and FLOWDROID threads. Note that the test machine is a shared system, but we made sure there are always enough resources for our evaluation available. Still, background services might influence the performance of a single run. To stamp out this factor, we ran each app three times. If there were outliers, we repeated the runs.

For this evaluation, we chose to use a non-default configuration of FLOWDROID. First, we disabled static field tracking due to the global scope as described in section 3.2. Next, instead of the EasyTaintWrapper, we use the SummaryTaintWrapper, which utilizes STUBDROID. STUBDROID’s precomputed data flow summaries are more precise than EasyTaintWrapper’s simple rules. We set the timeout for the data flow analysis to 10 minutes. The configuration summary is in Table 6.3.

We did not use the full sources and sinks list included in FLOWDROID because such would result in hundreds of sources and sinks per app and probably a long runtime. Instead, we chose to analyze which sensitive and possibly user-identifying data is sent out to the internet. As we want to compare the forwards and backward implementation, it is also

Option	Value
Array Size Tainting	disabled
Inspect Sources & Sinks	disabled
Static Field Tracking	disabled
Ignore Flows in System Packages	enabled
Exclude Soot Library Classes	enabled
Timeout	10 minutes
Taint Wrapper	SummaryTaintWrapper

Table 6.3.: Real World Apps Configuration

essential to not put one at a disadvantage. We opted for a 2:1 ratio of sources to sinks. This decision is based on the results of SuSI, a tool to automatically find sources and sinks in the Android framework [14]. Their extracted list of sources and sinks contains roughly 2.17 times more sources than sinks. The list of sources and sinks used in this evaluation is in Table 6.4 and Table 6.5.

We used FlowDroid’s forwards implementation on the to that date latest upstream commit<sup>2</sup> from the develop branch for the point of comparison. The backward implementation ran on our latest commit with all changes merged from the upstream.

We chose 200 apps randomly out of a Google Playstore dump from 2021 containing over 6000 apps for our evaluation set. The full list is appended to this work in ??.

### 6.2.2. Results

Out of 200 apps, 60 apps do not have any sources or sinks and thus, the analysis did not start. For 15 apps, the analysis aborted with errors. We are left with 125 apps for which both implementations completed the analysis.

Preliminary!  
Without  
memwatcher

### 6.2.3. Comparison to forwards analysis

Basically the answer to RQ1: Is the backwards search efficient enough to perform analysis on real world apps?

<sup>2</sup>The latest upstream commit was at that time b436733fc4a5130dfe4ce8ddb3f76fd374e9a487.

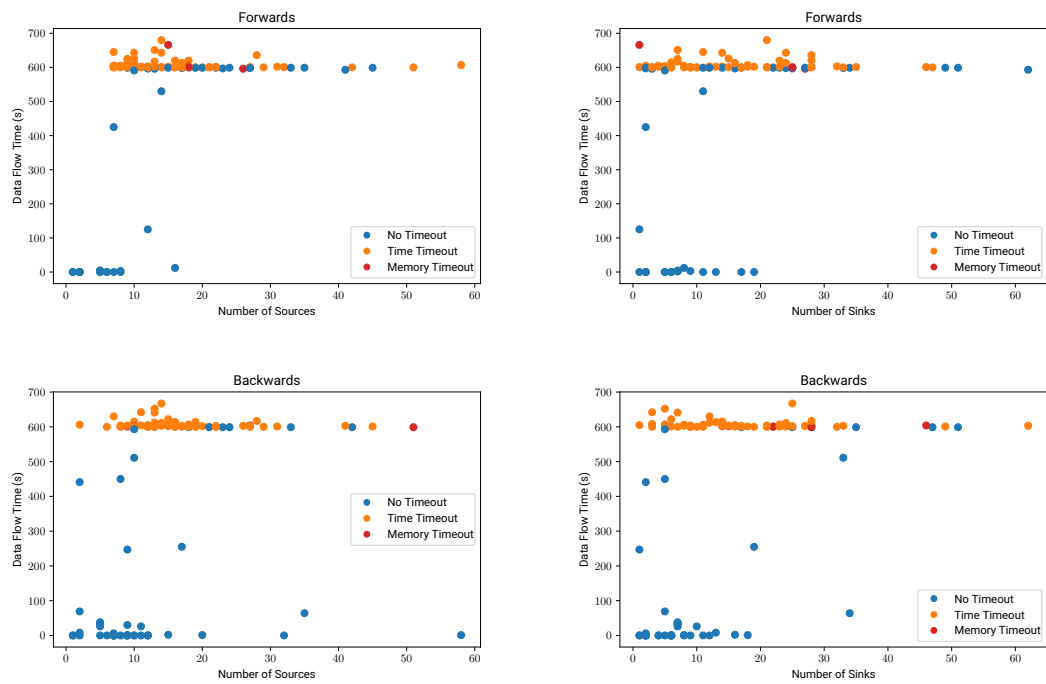


Figure 6.1.: Data flow time in comparison to source and sink count

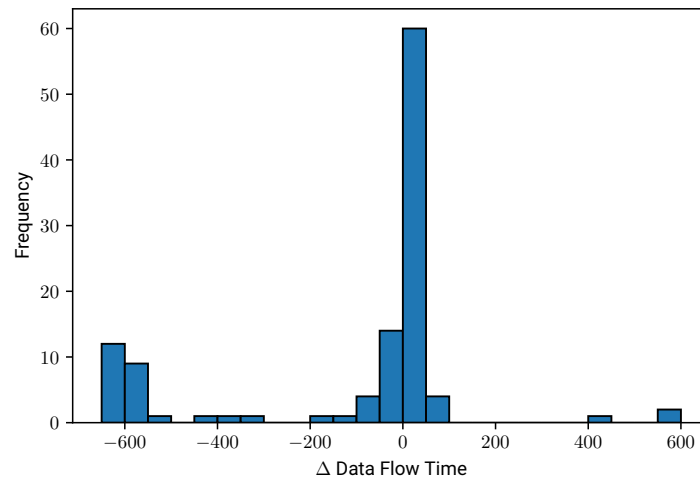


Figure 6.2.: Histogram of the Delta Data Flow Time

Basically the answer to RQ2: Can we find a pre-analysis known parameter to decide which analysis is more efficient?

Class	Method
android.location.Location	getLatitude() getLongitude()
android.location.LocationManager	getLastKnownLocation()
android.telephony.TelephonyManager	getDeviceId() getSubscriberId() getSimSerialNumber() getLine1Number() getImei() getMeid()
android.bluetooth.BluetoothAdapter android.net.wifi.WifiInfo	getAddress() getMacAddress() getSSID() getIpAddress()
java.net.InetAddress	getHostAddress()
android.telephony.gsm.GsmCellLocation	getCid() getLac()
android.content.pm.PackageManager	getInstalledApplications() getInstalledPackages() queryIntentActivities() queryIntentServices() queryBroadcastReceivers()
android.content.SharedPreferences android.provider.Browser	getDefaultSharedPreferences() getAllBookmarks() getAllVisitedUrls

Table 6.4.: Sources for Real World Apps Evaluation



Class	Method
java.net.URL	set() openConnection()
java.net.URLConnection	connect() setRequestProperty()
android.net.http.HttpsURLConnection	openConnection()
android.net.http.Headers	setEtag() setContentTypes() setLastModified() setLocation()
android.net.http.AndroidHttpClientConnection	sendRequestHeader()
android.net.http.RequestQueue	queueRequest()

Table 6.5.: Sinks for Real World Apps Evaluation

Metric	Forwards	Backwards
Average Runtime	504.7s	409.5s
Median Runtime	600.0s	600.0s
Average Abstractions Inflow	34651810	13741637
Average Abstractions Alias	12492775	33686397
Average Total Abstractions	47144585	47428034
Memory Timeouts	3.48%	4.35%
Time Timeouts	59.13%	53.91%

Table 6.6.: Results

	Forwards	Backwards
Average Runtime	331.6s	133.96s
Median Runtime	591.0s	1.0s
Average Abstractions Inflow	7468347	2262268
Average Abstractions Alias	2658090	4884307
Average Total Abstractions	10126437	7146575

Table 6.7.: Results without Timeouts

---

## 7. Related Work

---

We found very little about backward-directed taint analysis in our literature research.

Lerch et al.[11] contributed FlowTwist, a static taint analysis tool based on IFDS to detect confused deputy problems<sup>1</sup> in libraries. They identify the cause of such as a combination of an integrity and a confidentiality problem. For the integrity part, the sinks perform sensitive operations and the sources are attacker controlled while for the confidentiality part, the sinks can be read by an attacker and sources provide sensitive data. A combination of both naturally gives a centered statement. Now, the integrity sources and confidentiality sinks are way more frequent, thus they propose to solve the integrity part backwards and the confidentiality part forwards. In contrast to FLOWDROID, FlowTwist focuses on a special case of taint analysis and the applicability is quite narrow.

Yan et al.[21] proposed a vulnerability detection tool for PHP with a focus on web applications. They aim to detect typical web application vulnerabilities such as cross-site-scripting and SQL injections using backwards taint analysis. Instead of relying on reducing the problem to proven frameworks such as IFDS or IDE, they seemingly define their own data flow algorithm. The proposed algorithm traverses the basic blocks backwards and copies taints left after traversing the basic block to its predecessors. Unlike in our work and in general in data flow analysis, they do not try to reach a fixpoint, instead they just do not follow circular paths in the control-flow graph. They also emphasize their concept of "cleans": a predefined list of sanitization methods which kill the incoming taints. In FLOWDROID the same is possible using taint wrappers and both shipped implementations support such a concept. A rationale for traversing backwards, which is why we included it as related work, is not provided.

IFDS is not the only way to formulate a taint analysis.

---

<sup>1</sup>A confused deputy is a legitimate program with more privileges tricked into misusing its authority by a malicious program. Source: [https://en.wikipedia.org/wiki/Confused\\_deputy\\_problem](https://en.wikipedia.org/wiki/Confused_deputy_problem)

---

Synchronized pushdown systems (SPDS) by Späth et al.[18] are an alternative to IFDS with access paths for modelling a precise context-, flow- and field-sensitive data flow analysis. Similiar to IFDS it constructs a context-free grammar representing the call stack to ensure context-sensitivity. In addition, it also models field-sensitivity using a context-free grammar. Then it computes the acceptance state of both pushdown automaton to combine context- and field-sensitivity. This allows to represent recursive access paths such as `lst.next.prev.next...` without loss of precision where with IFDS and access paths,  $k$ -limiting is needed to ensure termination. SPDS are still an overapproximation in the case where at a statement both automaton are in acceptance state but via a different paths.

Doop [5] is another framework for points-to and, using P/Taint [7], also taint analysis. In contrast to others, it uses a declarative approach. Doop's frontend depends Soot to create facts and encodes them in tables. The analyses are encoded in a declarative rule set written in Datalog. These rule sets are then fed into the datalog solver Soufflé<sup>2</sup>.

---

<sup>2</sup><https://souffle-lang.github.io/>



---

## 8. Conclusion

---

---

## Bibliography

---

- [1] Alfred V. Aho, Ravi Sethi, and Jeffrey D. Ullman. *Compilers, principles, techniques, and tools*. eng. Reading, Mass. : Addison-Wesley Pub. Co., 1986. ISBN: 9780201100884. URL: <http://archive.org/details/compilersprincip00ahoa> (visited on 02/15/2021).
- [2] Steven Arzt. “Static Data Flow Analysis for Android Applications”. en. PhD thesis. Darmstadt: Technische Universität, 2017. URL: <https://tuprints.ulb.tu-darmstadt.de/5937/> (visited on 01/28/2021).
- [3] Steven Arzt et al. “FlowDroid”. In: *ACM SIGPLAN Notices* 49.6 (June 2014), pp. 259–269. DOI: 10.1145/2666356.2594299.
- [4] Eric Bodden. “Inter-procedural data-flow analysis with IFDS/IDE and Soot”. In: *Proceedings of the ACM SIGPLAN International Workshop on State of the Art in Java Program analysis - SOAP ’12*. ACM Press, 2012. DOI: 10.1145/2259051.2259052.
- [5] M. Bravenboer and Yannis Smaragdakis. “Strictly Declarative Specification of Sophisticated Points-to Analyses”. In: vol. 44. Oct. 2009, pp. 243–262. DOI: 10.1145/1640089.1640108.
- [6] A. Deutsch. “Interprocedural may-alias analysis for pointers: beyond k-limiting”. In: *PLDI ’94*. 1994. DOI: 10.1145/178243.178263.
- [7] Neville Grech and Yannis Smaragdakis. “P/Taint: unified points-to and taint analysis”. In: *Proceedings of the ACM on Programming Languages* 1 (Oct. 2017), pp. 1–28. DOI: 10.1145/3133926.
- [8] Neil Jones and Steven Muchnick. “Flow Analysis and Optimization of Lisp-Like Structures.” In: Jan. 1979, pp. 244–256. DOI: 10.1145/567752.567776.
- [9] Uday Khedker, Amitabha Sanyal, and Bhaurao Sathe. *Data Flow Analysis: Theory and Practice*. Jan. 2009. ISBN: 9780849332517. DOI: 10.1201/9780849332517.
- [10] Patrick Lam et al. “The Soot framework for Java program analysis: a retrospective”. en. In: Oct. 2011. URL: <http://www.bodden.de/pubs/lblh11soot.pdf> (visited on 03/11/2021).

- 
- 
- [11] J. Lerch et al. “FlowTwist: efficient context-sensitive inside-out taint analysis for large codebases”. In: *SIGSOFT FSE* (2014). DOI: 10.1145/2635868.2635878.
  - [12] Johannes Lerch and Ben Hermann. “Design your analysis: a case study on implementation reusability of data-flow functions”. In: *Proceedings of the 4th ACM SIGPLAN International Workshop on State Of the Art in Program Analysis*. ACM, June 2015. DOI: 10.1145/2771284.2771289.
  - [13] Nomair A. Naeem, Ondřej Lhoták, and Jonathan Rodriguez. “Practical Extensions to the IFDS Algorithm”. In: *Lecture Notes in Computer Science*. Springer Berlin Heidelberg, 2010, pp. 124–144. DOI: 10.1007/978-3-642-11970-5\_8.
  - [14] Siegfried Rasthofer, Steven Arzt, and E. Bodden. “A Machine-learning Approach for Classifying and Categorizing Android Sources and Sinks”. In: *NDSS*. 2014. DOI: 10.14722/NDSS.2014.23039.
  - [15] Thomas Reps, Susan Horwitz, and Mooly Sagiv. “Precise interprocedural dataflow analysis via graph reachability”. In: *Proceedings of the 22nd ACM SIGPLAN-SIGACT symposium on Principles of programming languages - POPL ’95*. ACM Press, 1995. DOI: 10.1145/199448.199462.
  - [16] H. Rice. “Classes of recursively enumerable sets and their decision problems”. In: (1953). DOI: 10.1090/S0002-9947-1953-0053041-6.
  - [17] Jonathan Rodriguez and Ondřej Lhoták. “Actor-Based Parallel Dataflow Analysis”. en. In: *Compiler Construction*. Ed. by Jens Knoop. Lecture Notes in Computer Science. Berlin, Heidelberg: Springer, 2011, pp. 179–197. ISBN: 9783642198618. DOI: 10.1007/978-3-642-19861-8\_11.
  - [18] Johannes Späth, Karim Ali, and Eric Bodden. “Context-, flow-, and field-sensitive data-flow analysis using synchronized Pushdown systems”. In: *Proceedings of the ACM on Programming Languages* 3 (Jan. 2019), pp. 1–29. DOI: 10.1145/3290361.
  - [19] Douglas Thain. *Introduction to Compilers and Language Design*. en. Google-Books-ID: 5mVyDwAAQBAJ. Lulu.com, July 2019. ISBN: 9780359138043.
  - [20] Raja Vallee-rai and Laurie Hendren. “Jimple: Simplifying Java Bytecode for Analyses and Transformations”. In: (Jan. 2004).
  - [21] X. Yan, H. Ma, and Q. Wang. “A static backward taint data analysis method for detecting web application vulnerabilities”. In: *2017 IEEE 9th International Conference on Communication Software and Networks (ICCSN)*. ISSN: 2472-8489. May 2017, pp. 1138–1141. DOI: 10.1109/ICCSN.2017.8230288.



---

## A. Appendix

---