

# Taint Analyse für Android Apps

Thomas Czogalik

Betreuer: Simon Greiner

## 1 Motivation

Durch fehlerhafte oder absichtlich manipulierte Programme können Datenlecks entstehen und diese ausgenutzt werden. In Abbildung 1 sehen wir ein Ausschnitt aus einem Java Programm. Es wird zunächst eine Datenbank Verbindung hergestellt und im Anschluß soll ein SQL Statement ausgeführt werden. Wie wir in Zeile 6 sehen, hängt das Statement von einer Benutzereingabe aus Zeile 5 ab. Ein Angreifer könnte in diesem Fall durch folgende Eingabe: "foo; DROP TABLE users" die Datenbank users löschen. Er könnte aber auch mit einem SELECT Statement Daten aus der Datenbank holen. Solche Szenarien lassen sich durch Überprüfung vermeiden. So ein Fehler ist aber nicht immer so einfach zu erkennen und wird übersehen. Da Computersysteme heutzutage in nahezu allen Bereichen unseres Lebens integriert sind, können solche Datenlecks fatale Folgen haben. Besonders Smartphones verwalten und verarbeiten viele vertrauliche und private Daten und kommunizieren dabei meist mit der Außenwelt. Im Februar 2015 befanden sich im Google Play Store ca. 1.4 Millionen Apps. Diese sind jedem zugänglich, der auf seinem mobilen Gerät das Betriebssystem Android installiert hat. Bei so einer großen Anzahl Apps bietet der Google Play Store eine große Angriffsfläche. Deshalb ist es notwendig, den Fluss sensibler Daten nachvollziehen zu können. Dies ist mithilfe einer Taint Analyse möglich. Der Autor möchte im folgenden zunächst auf die Grundlagen der Taint Analyse eingehen und danach die Probleme vorstellen, die bei einer Taint Analyse für Android Apps auftreten. Im Anschluß wird FlowDroid vorgestellt eine Taint Analyse für Android Apps.

---

```
1 connection = ...
2 stmt = connection.createStatement();
3 BufferedReader br = ...
4 String name = br.readLine();
5 String sql = "SELECT * FROM users WHERE name=" + name + ";";
6 stmt.executeQuery(sql);
```

---

**Fig. 1.** Beispiel bei dem SQL Injections möglich sind

## 2 Taint Analyse

### 2.1 statisch vs. dynamisch

Bei der Taint Analyse unterscheidet man zwischen statischer und dynamischer Taint Analyse. Der Vorteil der statischen Variante ist, dass das Kompilat nicht ausgeführt werden muss. Dies ist Hilfreich, da heutige Malware erkennen kann ob sie überwacht wird und kann ihr Verhalten anpassen. Im folgenden wird sich auf die statische Taint Analyse beschränkt.

### 2.2 Spezifikation und Vorgehen

Die Idee der Taint Analyse ist, dass jede von außen veränderbare Variable ein Sicherheitsrisiko birgt. Ihr Ziel ist es die Software gegen externe Angriffe sowie interne Risiken abzusichern. Dazu sucht die Taint Analyse nach Datenflüssen von möglichen tainted sources zu einem sink. Als source wird eine Funktion bezeichnet, die Quelle sensibler Daten ist. Ein sink ist eine Funktion, die Daten möglicherweise an nicht vertrauenswürdige Beobachter weitergibt. Welche Funktionen in einzelnen sources und sinks sind, muss vor der Taint Analyse angegeben werden.

### 2.3 Formal

Im Folgenden wird eine Formalisierung der Taint Analyse vorgestellt mithilfe von Schlussregeln (2). Die SOURCE Regel beschreibt die Einfuhr von tainted Daten durch einen source.  $src(m, l)$  bedeutet, dass dem Parameter  $m$  Daten aus einer source zugewiesen werden mit dem Namen  $l$ .  $l$  stammt aus einem Source-Label-Set. Wenn nun auch  $m$  auf einen Objektrepräsentanten  $o$  zeigt, ist  $o$  durch die Source  $l$  tainted. Als Objektrepräsentant wird ein Objekt im Sinne der Analyse bezeichnet. Da nicht sicher ist auf wieviele Objekte tatsächlich zur Laufzeit ein Pointer zeigt, fasst man diese Menge zusammen und wählt einen Objektrepräsentanten. Die zweite Regel TRANSFER beschreibt die Übertragung von tainted Daten.  $transfer(m, n)$  bedeutet, dass Daten vom Parameter  $m$  zu  $n$  übertragen werden. Wenn außerdem der Objektrepräsentant  $o1$  von einer Source  $l$  tainted ist und  $m$  auf  $o1$  und  $n$  auf ein  $o2$  zeigt, dann ist auch  $o2$  durch  $l$  tainted. Die letzte Regel SINK beschreibt den Datenfluss von verschmutzten Daten zu einem Sink. Wenn der Objektrepräsentant  $o$  durch eine Source  $so$  tainted ist und die Variable  $m$ , die auf  $o$  zeigt einem Sink mit dem Namen  $si$  übergeben wird, dann gibt es einen Datenfluss von der Source  $so$  zum Sink  $si$ .  $si$  stammt aus einem Sink-Label-Set. Mit diesen drei Regeln kann man eine Taint Analyse durchführen indem man diese Fixpunkt Iterativ ausführt, bis nichts mehr gefunden wird.

**2.3.1 Beispiel** Wendet man diese Regeln an dem Beispiel aus Abbildung 1 an, bekommt man folgendes Ergebnis. Zunächst müssen die Sinks und Sources

$$\begin{array}{c}
\frac{src(m, l) \quad m \rightarrow o}{tainted(o, l)} \text{ (SOURCE)} \\
\\
\frac{tainted(o_1, l) \quad m \rightarrow o_1 \quad n \rightarrow o_2 \quad transfer(m, n)}{tainted(o_2, l)} \text{ (TRANSFER)} \\
\\
\frac{tainted(o, so) \quad m \rightarrow o \quad sink(m, si)}{flow(so, si)} \text{ (SOURCE)}
\end{array}$$

**Fig. 2.** Schlussregeln für die Taint Analyse

definiert werden. Dazu wird die Funktion *readLine* aus Zeile 4 in das Source-Label-Set eingefügt. In das Sink-Label-Set wird die Funktion *executeQuery* eingefügt. Die Taint Analyse wird zunächst die SOURCE Regel in Zeile 4 anwenden, da dort Daten aus einer Source dem String *name* zugewiesen werden. Somit wird der String *name* als tainted markiert. Als nächstes wird die TRANSFER Regel in Zeile 5 angewendet. Hier werden Daten vom String "name" zum String *sql* übertragen. Deshalb wird auch hier der String *sql* als tainted markiert. In Zeile 5 wird die SINK Regel angewendet. Hier wird der String *sql* in den Sink *executeQuery* übergeben. Da *sql* tainted ist, wird der Sink *executeQuery* auch als tainted markiert. Somit wurde ein tainted Datenfluss von der Source *readLine* zum Sink *executeQuery* gefunden.

## 2.4 Sanitization

Nachdem die Taint Analyse Datenflüsse von sources zu sinks entdeckt hat, sollte dem Programmierer die Möglichkeit gegeben werden, seine eingehenden Daten zu überprüfen und zu säubern. In dem Beispiel aus Abbildung 1 würde zum Beispiel die Funktion *onlyLetters(s : String) : String* die Benutzereingabe nehmen und alle nicht Buchstaben aus der Eingabe entfernen und einen String zurückgeben, der nur aus Buchstaben besteht. Einige Programmiersprachen bieten solche Säuberungsfunktionen an. PHP hat zum Beispiel die Funktion *htmlentities*, die Zeichen, die in HTML besondere Bedeutung haben in ihre HTML Instanzen konvertiert. Aus dem Zeichen '<' wird '&lt;'. Neben den eingebauten sanitization Funktionen können auch eigene Funktionen geschrieben werden. Dies kann jedoch gefährlich werden, da Daten die einen säuberungsprozess durchlaufen haben als harmlos angesehen werden. Das kann zu Problemen führen, wenn die Funktion nicht vollständig oder falsch ist, da es im vornherein schwierig festzulegen ist welche Fälle durch eine Sanitization Funktion beachtet werden müssen. Außerdem müssen die Sanitization und die sink Funktion zusammen passen. In Abbildung 1 würde die Funktion *htmlentities* die Eingabe nicht säubern.

## 2.5 Präzision

Eine gute Analyse sollte in der Lage sein false positives nicht zu melden. Dadurch wird die manuelle aussortierung dieser durch den Benutzer vermieden. Als false positives werden Fehlermeldungen bezeichnet, die einen Fehler anzeigen, obwohl keiner vorliegt. Mit den folgenden Eigenschaften lassen sich false positives vermeiden.

**2.5.1 Fluss-Sensitivität** In Abbildung 3 wird einem String  $s$  in Zeile 2 zunächst ein harmloser Wert zugewiesen. Im Anschluss wird dieser String einem Sink übergeben. Obwohl die tainted Daten, die dem String  $s$  in Zeile 4 zugewiesen werden nie in den sink kommen, wird dieser als tainted markiert. Grund dafür ist, dass die Reihenfolge nicht beachtet wird. Beachtet man diese findet die Analyse den false positive nicht. Diese Eigenschaft nennt man Fluss-Sensitivität.

**2.5.2 Kontext-Sensitivität** Obwohl  $s_2$  in Abbildung 4 keine tainted Daten zugewiesen bekommt, wird der sink in Zeile 4 als tainted markiert. Grund dafür ist die Funktion  $id(s : String) : String$  aus Zeile 7. Da  $s_1$  und  $s_2$  die Funktion  $id$  aufrufen, zeigt das  $s$  aus Zeile 8 sowohl auf  $s_1$  als auf  $s_2$ . Da  $id$  mit den tainted Daten aus dem Source aufgerufen wird, werden  $s_1$  und  $s_2$  als tainted markiert. Durch Beachtung des Aufrufkontexts, zeigt das  $s$  aus Zeile 8 entweder auf  $s_1$  oder auf  $s_2$ . Somit wird bei dieser herangehensweise kein Fehler gefunden.

**2.5.3 Objekt-Sensitivität** Ein ähnliches Problem wie zuvor ist in Abbildung 5 zu sehen. Die Funktion  $getValue() : String$  der Klasse Foo zeigt sowohl auf  $o_1$  als auf  $o_2$ . Obwohl der Sink in Zeile 7 keine tainted Daten erhält wird hier ein Fehler angezeigt, da nicht zwischen den einzelnen Objekten unterschieden werden kann. Eine Analyse die Objekt-Sensitiv ist findet solche false positives nicht. Die Objekt-Sensitivität kann auch als Subtyp der Kontext-Sensitivität gesehen werden.

**2.5.4 Feld-Sensitivität** In der letzten Abbildung (6) ist ein Objekt  $o$ , dass mehrere Felder hat. Einem der Felder werden tainted Daten aus einer Source zugewiesen. Dem anderen Feld werden harmlose Daten zugewiesen. Obwohl das harmlose Feld dem Sink übergeben wird, wird der Sink als tainted markiert. Der Grund ist, dass Felder bei Feld-Insensitiven Analysen zu ihrem Basis Objekt zusammengeführt werden und als ganzes betrachtet werden. Da hier  $field_2$  tainted ist, wird das ganze Objekt  $o$  tainted und somit auch das  $field_1$ . Bei einer Feld-Sensitiven Analyse werden die Felder einzeln betrachtet und ein Fehler wird hier nicht gefunden.

---

```
1 void flowSensitive() {  
2     String s = "value";  
3     sink(s);  
4     s = source();  
5 }
```

---

**Fig. 3.** Fluss-Sensitivität

---

```
1 void contextSensitive() {  
2     String s1 = id(source());  
3     String s2 = id("123");  
4     sink(s2);  
5 }  
6  
7 String id(String s) {  
8     return s;  
9 }
```

---

**Fig. 4.** Kontext-Sensitivität

---

```
1 void objectSensitive() {  
2     Foo o1 = new Foo();  
3     Foo o2 = new Foo();  
4     o1.value = "123";  
5     o2.value = source();  
6  
7     sink(o1.getValue());  
8 }  
9  
10 class Foo {  
11     String value;  
12  
13     String getValue() {  
14         return value;  
15     }  
16 }
```

---

**Fig. 5.** Objekt-Sensitivität

---

```

1 void fieldSensitive() {
2     Foo o = new Foo();
3     o.field1 = "123";
4     o.field2 = source();
5
6     sink(o.field1);
7 }
8
9 class Foo {
10     String field1;
11     String field2;
12 }

```

---

**Fig. 6.** Feld-Sensitivität

### 3 Android

Android ist sowohl ein Betriebssystem als auch eine Software-Plattform für mobile Geräte. Es bietet den Entwicklern eine Schnittstelle an, um auf verschiedene System Funktionalitäten des mobilen Gerätes zuzugreifen. Der Nutzerstandort ist zum Beispiel eine solche Funktionalität. Der Entwickler muss die von ihm genutzte Funktionalität in die *AndroidManifest.xml* eintragen. Der Benutzer, der die App installieren möchte muss dieser die Berechtigung geben diese Funktionalitäten verwenden zu dürfen. Werden die Berechtigungen verweigert, kann die App nicht installiert werden. Das Problem ist, dass oft zu viele Berechtigungen gefordert werden. Entweder aus Unkenntnis oder mit bösen Absichten. Außerdem sind viele Berechtigungen zu mächtig und können ausgenutzt werden. Eine social-network App zum Beispiel, möchte die Berechtigung *android.permission.READ.CONTACTS* um neue Freunde anhand der Email Adresse vorschlagen zu können. Sie kann aber mit dieser Berechtigung auch die Telefonnummern und sonstige Kontaktdaten auslesen, auch wenn der Kontakt nicht bei dem Dienst angemeldet ist. Trojaner können dies ausnutzen. Um dem entgegenzuwirken kann man die Datenflüsse einer Android App mithilfe der Taint Analyse analysieren. Dabei gibt es einige Besonderheiten und Schwierigkeiten gegenüber anderen *JAVA* Programmen. Eine Android App kann aus vier verschiedenen Komponenten bestehen.

Activity	Einzelner Screen, der für den Benutzer sichtbar ist
Service	Eine Aktion die im Hintergrund abläuft
Content Provider	Ist zuständig für das Speichern von Daten
Broadcast Reciever	Wartet auf ein globales Ereignis und führt eine vodefinierte Aktion durch

Aufgrund dieser Besonderheit entstehen bestimmte Probleme für eine Taint Analyse. Zum einen hat eine Android App keine zentrale main Methode son-

dern mehrere Einstiegspunkte. Eine App kann zum Beispiel aus drei Activities und einem Service bestehen. Es gibt zwar eine Haupt-Activity, aber welche Komponente als nächstes an die Reihe kommt ist unbestimmt, da es meißt auf die Benutzereingabe ankommt. Außerdem haben die einzelnen Komponenten besonderes Verhalten. Sie können zum Beispiel gestartet und beendet werden oder angehalten werden, wenn der Speicher voll ist und fortgesetzt, wenn wieder Speicher frei ist. Somit entsteht ein komplizierter Lebenszyklus, mit dem die Taint Analyse umgehen muss. Außerdem muss vor der Analyse spezifiziert werden, was sinks und sources in Android sind.

## 4 FlowDroid

Diesen Problemen haben sich die Entwickler von FlowDroid gestellt und haben eine statische Taint Analyse die, Fluss-, Kontext-, Objekt- und Feld-Sensitiv ist implementiert <sup>1</sup>. Für die Sinks und Sources benutzt FlowDroid die Ausgabe des Tools SuSi. SuSi ist ein machine-learning Tool, dass vollautomatisch den Android Source Code analysiert und eine Liste von Sinks und Sources ausgibt. Außerdem wird ein präziser Android Lebenszyklus modelliert.

### 4.1 Lebenszyklus

Aufgrund der vielen möglichen Einstiegspunkte bei einer Android App, wird bei FlowDroid zunächst eine dummy-main-Methode erstellt. Diese Methode wird individuell für jede App erstellt. Die dummy-main enthält nur Teile des Lebenszyklus die auch auftreten können. Um diese zu indentifizieren werden die XML Konfigurationsdateien und der Quellcode analysiert. Danach ruft die dummy-main-Methode nacheinander Teile des Zykluses auf.

### 4.2 Evaluation

Für die Evaluierung wurde die Testumgebung DroidBench <sup>2</sup> verwendet. DroidBench enthielt zum damaligen Zeitpunkt 39 Apps mit verschiedenen Analyse Problemen. Verglichen wurde FlowDroid mit den beiden kommerziellen Tool App Scan von IBM und Fortify SCA von HP. In dem Evaluierungsbogen (7) ist zu erkennen, dass FlowDroid in den Kategorien Callbacks und Lifecycle deutlich besser ist, als die beiden kommerziellen Tools. Dies ist auf die präzise Modellierung des Lebenszyklus zurückzuführen. In ein paar Tests schlägt FlowDroid aber auch fehl. Zum Beispiel findet FlowDroid keinen Fehler im IntentSink1 Test. Das liegt daran, dass in dem Test tainted Daten aus der einen Komponente in die nächste übergeben wird und diese die Daten direkt wieder zurück schickt. FlowDroid verfolgt tainted Daten nur innerhalb der Komponenten. Deshalb schlägt

<sup>1</sup> <https://github.com/secure-software-engineering/soot-infoflow-android/wiki>

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

der Test hier fehl. Ein weiterer Test der Fehlschlägt ist der StaticInitialization1 Test. Hier liegt der Grund in dem verwendeten Framework Soot. Soot ist ein Framework zum analysieren und transformieren von Java und Android Programmen. Es kann aus Java Code einen Call-Graphen generieren. Dabei geht Soot davon aus, dass die komplette statische Initialisierung zu beginn stattfindet, was in diesem Testfall nicht der Fall ist. Im Ergebniss schneidet FlowDroid jedoch besser ab, als seine beiden Konkurrenten. FlowDroid ist mit 86% etwas besser beim Verhältnis zwischen gefunden Fehlern und false-positives (Precision p). Mit 93% gegenüber 50% und 61% ist FlowDroid aber deutlich besser beim auffinden aller Fehler (Recall r). FlowDroid testete zusätzlich die damals Top 500 Apps des Google Play Stores und fand in den meisten Apps ein Datenleck. Es waren jedoch keine bösen Absichten erkennbar. Außerdem wurden 1000 Malware Apps des VirusShare Projects analysiert. Es wurden pro App ca. zwei Datenlecks gefunden.

## 5 Ausblick

FlowDroid macht schon viele richtig, aber wie die Evaluierung gezeigt hat, ist FlowDroid noch nicht perfekt. Zum Einen sollte die Inter-App-Kommunikation verbessert werden, damit auch tainted Daten zwischen den Komponenten verfolgt werden können. Dazu könnte der Sink der einen Komponente gleichzeitig die Source der anderen Komponenten sein. Zum Anderen hat FlowDroid keine Sanitization implementiert. Dies ist jedoch notwendig, damit der Programmierer auf gefundene Fehler reagieren kann. Für die Zukunft wäre es auch wichtig, dass FlowDroid unterscheiden kann zwischen notwendigen Datenlecks und böartigen. Eine Navigationsapp benötigt den Benutzerstandort. Eine Taschenlampen App nicht. Schließlich wäre es auch hilfreich, wenn die Analyse für bestimmte Arten von Sinks und Sources einschränkbar wäre. Der [?]Vorteil wäre, dass man sich bestimmte Datenflüsse näher anschauen könnte. SuSi klassifiziert und kategorisiert bereits die Sinks und Sources.

## References



⊕ = correct warning, ★ = false warning, ○ = missed leak  
multiple circles in one row: multiple leaks expected  
all-empty row: no leaks expected, none reported

App Name	AppScan Source	Fortify SCA	FlowDroid
<b>Arrays and Lists</b>			
ArrayAccess1			★
ArrayAccess2	★	★	★
ListAccess1	★	★	★
<b>Callbacks</b>			
AnonymousClass1	○	⊕	⊕
Button1	○	⊕	⊕
Button2	⊕ ○ ○	⊕ ○ ○	⊕ ⊕ ⊕ ★
LocationLeak1	○ ○	○ ○	⊕ ⊕
LocationLeak2	○ ○	○ ○	⊕ ⊕
MethodOverride1	⊕	⊕	⊕
<b>Field and Object Sensitivity</b>			
FieldSensitivity1			
FieldSensitivity2			
FieldSensitivity3	⊕	⊕	⊕
FieldSensitivity4	★		
InheritedObjects1	⊕	⊕	⊕
ObjectSensitivity1			
ObjectSensitivity2	★		
<b>Inter-App Communication</b>			
IntentSink1	⊕	⊕	○
IntentSink2	⊕	⊕	⊕
ActivityCommunication1	⊕	⊕	⊕
<b>Lifecycle</b>			
BroadcastReceiverLifecycle1	⊕	⊕	⊕
ActivityLifecycle1	⊕	⊕	⊕
ActivityLifecycle2	○	⊕	⊕
ActivityLifecycle3	○	○	⊕
ActivityLifecycle4	○	⊕	⊕
ServiceLifecycle1	○	○	⊕
<b>General Java</b>			
Loop1	⊕	○	⊕
Loop2	⊕	○	⊕
SourceCodeSpecific1	⊕	⊕	⊕
StaticInitialization1	○	⊕	○
UnreachableCode		★	
<b>Miscellaneous Android-Specific</b>			
PrivateDataLeak1	○	○	⊕
PrivateDataLeak2	⊕	⊕	⊕
DirectLeak1	⊕	⊕	⊕
InactiveActivity	★	★	
LogNoLeak			
<b>Sum, Precision and Recall</b>			
⊕, higher is better	14	17	26
★, lower is better	5	4	4
○, lower is better	14	11	2
Precision $p = \frac{\oplus}{(\oplus + \star)}$	74%	81%	86%
Recall $r = \frac{\oplus}{(\oplus + \circ)}$	50%	61%	93%
F-measure $2pr/(p+r)$	0.60	0.70	0.89

Fig. 7. DroidBench Evaluierung, aus[.]