

# Datacentric Semantics for Verification of Privacy Policy Compliance by Mobile Applications

Agostino Cortesi<sup>1,2</sup>, Pietro Ferrara<sup>1</sup>, Marco Pistoia<sup>1</sup>, and Omer Tripp<sup>1</sup>

<sup>1</sup> IBM Thomas J. Watson Research Center, USA

<sup>2</sup> Università Ca' Foscari Venezia, Italy

**Abstract.** We introduce an enhanced information-flow analysis for tracking the amount of confidential data that is possibly released to third parties by a mobile application. The main novelty of our solution is that it can explicitly keep track of the footprint of data sources in the expressions formed and manipulated by the program, as well as of transformations over them, yielding a lazy approach with finer granularity, which may reduce false positives with respect to state-of-the-art information-flow analyses.

**Keywords:** Abstract Interpretation, Privacy, Information-flow Analysis.

## 1 Introduction

Mobile applications typically ask for permission to access personal (i.e. relevant with respect to privacy) information stored on the device. However, even in non-malicious applications, once these permissions are granted it is often the case that data concerning gender, sex, age, GPS location, smartphone ID, etc. is managed in a way that partially releases it to third parties (e.g. for advertising, profiling, analytics and social computing), with or without some degree of obfuscation, leaving the user unaware of how much confidential information actually leaked [29,36]. Most systems, in fact, are designed to allow users to configure access control (e.g., by setting permissions), without enabling them to monitor the actual information flow of confidential data. In reality, users may trust an application to manage their personal information, but might be concerned about the obfuscation degree applied to that information before it is passed to other (possibly untrusted) actors. The key issue is to keep track of (and possibly restrict) the amount of confidential information that is released by an app, without compromising the usability of the app itself by enforcing overly conservative constraints.

In this challenging context, the aim of this work is to define a theoretical framework to support the design of tools that provide developers as well as end users with better control of how the values managed by the applications reveal confidential data stored on the device.

### 1.1 Background

Preserving confidentiality of sensitive information in software systems is a subject of intensive research. Various language-based information-flow security analyses were

proposed [18,21,32,35]. Most of these works are based on the non-interference notion that says that a variation of confidential data given as input to a program does not cause any variation of publicly observable data [12]. The approaches are different (type systems [28,32,34], dependence graphs [18,22], slicing [2,6,21,35], etc.), and apply to different languages including imperative, object-oriented, functional and structured query languages [17,18,21,30]. Recent works address in particular data protection from permission-hungry Android applications [4,7,15,19,24,40], and data-leakage aggregation due to undesired inter-application dataflows [31].

However, in the scenario depicted above the crucial point is not just to discover if sensitive data is confined to private variables, but also to keep control of how the authorized access to confidential data is compliant with respect to a privacy policy, expressed in terms of minimal degree of obfuscation that should be applied to sensitive data in the exposed values. In this respect, the information-flow approach of the mentioned works yields overly conservative results, as the granularity of public/private variables is too coarse, just like the tampered/untampered granularity [1] assigned to data when declassification mechanisms are introduced for relaxing confidentiality policies [5].

## 1.2 Contribution

This work extends taint analysis, which is a popular variant of information-flow analysis [37,38,39] in order to trace the dependence flow of confidential information from data sources to data sinks. A finer granularity of the analysis is obtained by explicitly keeping track of the footprint of data sources in the expressions managed by the program, as well as of the obfuscation impact of the program operators. As the analysis is defined as an instance of the Abstract Interpretation framework [10], the tradeoff between accuracy and efficiency can be tuned by a suitable choice of the concrete domain abstraction.

The main contributions of this paper can be summarized as follows:

- We design an enhanced concrete semantics that makes explicit the dependence of values on local data sources.
- We define the notion of “confidentiality value of an expression” in terms of min/max confidentiality degree of its sources and of min/max obfuscation degree of the operators that are used to generate it.
- We lift the enhanced concrete semantics to (computable) abstract semantics according to the Abstract Interpretation framework.
- We show how a static analysis based on this framework can be used to verify the satisfaction of privacy policies.
- We provide practical evidence of the effectiveness of our approach.

In concrete implementations, tracking indirect information flows negatively impacts the effectiveness of the analysis due to the presence of exceptions. Therefore, our approach disregards them. However, the treatment of implicit flows can be further incorporated into our framework by infusing relational operators’ footprint in the different conditional statements’ branches.

```

1 public class IMBanner {
2   public void loadBanner() {
3     UserInfo user = new UserInfo();
4     user.updateInfo();
5     BannerView banner = new BannerView(user);
6     banner.loadNewAd();
7     show(banner);
8   }
9 }
10
11 public class BannerView {
12   private UserInfo user;
13   BannerView(UserInfo user) {
14     this.user = user;
15   }
16   void loadNewAd() {
17     String url = "http://www.inmobi.com/...?id="
18       + user.id + "&lang="+user.language+
19       "&country=" + user.country + "&loc=" + user.loc;
20     //open an http connection with url
21     //update the new ad to display
22   }
23 }
24
25 public class UserInfo {
26   String language;
27   String country;
28   String id;
29   Location loc;
30
31   void updateInfo() {
32     Locale localLocale = Locale.getDefault();
33     language = localLocale.getLanguage();
34     country = localLocale.getCountry();
35     String androidId = Settings.Secure.getAndroidId();
36     id = MessageDigest.hashSHA1(androidId);
37     loc = LocationManager.getLastKnownLocation();
38   }
39 }

```

Fig. 1. Code Snippet from the Inmobi Library

### 1.3 Structure of the Paper

The rest of paper is structured as follows: Section 2 presents some examples that motivate the main novelties of our approach. Sections 3, 4, and 5 describe the syntax, and the enhanced concrete and abstract semantics, respectively. Section 6 introduces the notion of confidentiality and obfuscation values for sources and operators. Section 7 shows how this framework can be applied to the verification of privacy compliance policies. Section 8 discusses related work, while Section 9 concludes.

## 2 Motivating Examples

### 2.1 Inmobi

Consider the motivating example in Figure 1. This code is extracted from the *Inmobi* library.<sup>1</sup> *Inmobi* is among the three most popular advertisement engines for Android apps [3]. This code sketches the main steps performed by the *Inmobi* library when loading an advertisement banner. This is performed by method `IMBanner.loadBanner()`, that first creates and updates a `UserInfo` object (lines 3-4), then creates and loads the advertisement banner view passing the information about the user (lines 5-6), and finally displays the banner (line 7). Even though, at a first glance, this method does not seem to access any confidential information, `user.updateInfo()` collects and transforms various pieces of information about the user and the device, and in particular (i)

<sup>1</sup> This library is obfuscated, and some parts (and in particular `BannerView.loadNewAd()`) cannot be decompiled. For the sake of readability, we represented the main components of the library in this code snippet.

the language and country of the user from the default `Locale` object (lines 32-34), (ii) the hashing of the Android ID (lines 35-36)<sup>2</sup>, and (iii) the last known location (line 37). When a `UserInfo` object is passed to the constructor of `BannerView`, it is stored in a local field. The information contained in this field is then concatenated in the URL (lines 17-19) used to retrieve the advertisement banner. In this way, the data collected by `UserInfo.updateInfo()` is leaked to the advertisement server. This data aggregates various sources (language and country from the `Locale`, the android ID, and the location).

This example shows that we need to track complex flows of information. For instance, with a standard taint analysis [39] an alarm would be raised upon any flow from a source to a sink. In this particular example, we expect that it would be fine to release to the advertisement server some of the sources (e.g., the country and the language, but not the Android ID), so taint analysis could raise a false alarm in this scenario. Indeed, we are interested in computing the global amount of data that is released (that is, country, language, Android ID, and location), and to raise an alarm only if this amount exceeds a specified threshold. In addition, one needs to specify if the transformation performed on the confidential data (e.g., hashing the Android ID) is obfuscating the value sufficiently or not. For instance, the hashing of the ID might be used to track a user or device if the hash clashing is quite rare, and therefore the level of obfuscation performed by this transformation might be insufficient.

## 2.2 IMSI

The following code snippet is extracted from internal Android library `com.android.internal.telephony.cdma.RuimRecords`:

```
1 String mImsi = telephonyManager.getDeviceId();
2 log("IMSI:" + mImsi.substring(0, 6) + "xxxxxxx");
```

It leaks a portion of the device identifier through the log. The fundamental question here is whether the first 6 characters of the International Mobile Subscriber Identity (IMSI) code contain confidential information that the user does not want to leak outside. The IMSI code is usually made by 15 characters, where the first 3 characters identify the country, the following 2 or 3 characters identify the mobile network, and the rest is used to identify the device.<sup>3</sup> Therefore, we assume that the first 6 characters do not contain confidential information.

## 3 Syntax

At the lowest level of the language, we consider expressions on strings ( $s \in \mathcal{S}$ ), integers ( $n \in \mathbb{Z}$ ), and Boolean values ( $b \in \mathbb{B}$ ). We denote by *bexp* and *nexp* Boolean

<sup>2</sup> The Android ID is “randomly generated when the user first sets up the device and should remain constant for the lifetime of the user’s device”. Therefore, it is used to track a specific user (rather than a device) by advertisement engines. We simplify the API call to make the code more readable.

<sup>3</sup> [http://en.wikipedia.org/wiki/International\\_mobile\\_subscriber\\_identity](http://en.wikipedia.org/wiki/International_mobile_subscriber_identity).

and numerical expressions, respectively. In addition to basic numerical, textual, and Boolean expressions, we introduce label constants (that refer to datastore entries). Let  $\text{Lab}$  be the set of labels. We denote by  $\ell$ , possibly subscripted, labels identifiers in  $\text{Lab}$ . Our language support a statement  $\text{read}(\ell)$  that returns the value read from the datastore corresponding to the label  $\ell$ . We define string expressions by  $\text{sexp} ::= s \mid \text{sexp}_1 \circ \text{sexp}_2 \mid \text{encrypt}(\text{sexp}, k) \mid \text{sub}(\text{sexp}, n_{exp_1}, n_{exp_2}) \mid \text{hash}(\text{sexp}) \mid \text{read}(\ell_{exp})$  where  $\circ$  denotes the concatenation of two strings,  $k$  denotes a key used to encrypt a textual value,  $\text{sub}(s, n_1, n_2)$  computes the substring of  $s$  from the  $n_1$ -th to the  $n_2$ -th character, and  $\text{hash}(s)$  computes the hash value of  $s$ . For the sake of simplicity, we focus our formalization on this minimal representative language, and on operations over strings. Our approach can be extended straightforwardly to support other operations and types.

Finally, we define a standard minimal imperative set of statements. In particular, we support string assignment ( $x := \text{sexp}$ ), concatenation ( $c_1; c_2$ ), conditional if (if  $\text{bexp}$  then  $c_1$  else  $c_2$ ), and while loops (while  $\text{bexp}$  do  $c$ ). In addition, we have a special statement  $\text{send}(\text{sexp})$  that leaks a string value.

## 4 Collecting Semantics

### 4.1 Domain

First of all, we define atomic data expressions by  $\mathbb{D} = \{\langle \ell_i, L_i \rangle : i \in I\}$ . Given a set of data labels, which identify the locations of the *read-only*<sup>4</sup> datastore a program interacts with, an atomic data expression  $\text{adexp}$  is a set of elements  $\langle \ell_i, \{(op_j, \ell'_j) : j \in J\} \rangle$ . An element  $\langle \ell_i, \{(op_j, \ell'_j) : j \in J\} \rangle$  in  $\text{adexp}$  says that the value of  $\text{adexp}$  has been obtained from the datum stored in the location  $\ell_i$  by combining it with data coming from the locations  $\ell'_j$  through the corresponding operations  $op_j$ . In other words, an atomic data expression keeps track, for each source of the expression value, of the set of other data sources that were used to get that value from it. We denote by  $\mathbb{D}$  the domain of atomic data.

We focus our collecting semantics on the variables referring to values coming from the datastore. Therefore, we define a data environment mapping local variables in  $\text{Var}$  to atomic data expressions ( $D : \text{Var} \rightarrow \wp(\mathbb{D})$ ). Note that each variable may contain data about different sources (e.g., the concatenation of the strings representing the Android identifier and the location), and therefore each variable is related to a set of atomic data expressions. In addition, the concrete state tracks value information as well ( $V : \text{Var} \rightarrow (\mathbb{Z} \cup \mathbb{S})$ ). Formally,  $\Sigma = D \times V$ .

We then introduce a concrete datastore that contains all the possible atomic data that may be read by a program, where the special label  $\star$  is used to represent data coming either from the input of the program or from the constant set of the program itself, i.e. data that is not contained in the datastore.

**Definition 1 (Concrete Datastore).** *A concrete datastore  $C$  is a set  $\{\langle \ell_i, \emptyset \rangle : i \in I\} \subseteq \mathbb{D}$  such that  $\forall i, j \in I : i \neq j \Rightarrow \ell_i \neq \ell_j$ , and  $\ell_i \neq \star$ .*

<sup>4</sup> We restrict our focus to a read-only datastore, for the sake of simplicity. Extending the model to the general case brings about the problem of aliasing that should be studied further.

$$\begin{aligned}
S_A[x](a, v) &= a(x) \\
S_A[\text{read}(lexp)](a, v) &= \{ \langle S_L[\text{lexp}](a, v), \emptyset \rangle \} \\
S_A[\text{encrypt}(sexp, k)](a, v) &= \{ \langle \ell_1, L_1 \cup \{ \langle \text{encrypt}, k \rangle, \ell_1 \} \rangle : \langle \ell_1, L_1 \rangle \in S_A[\text{sexp}](a, s, n) \} \\
S_A[s](a, v) &= \{ \langle *, \emptyset \rangle \} \\
S_A[\text{sexp}_1 \circ \text{sexp}_2](a, v) &= \{ \langle \ell_1, L_1 \cup \{ \langle \circ, \ell_2 \rangle \} \rangle, \langle \ell_2, L_2 \cup \{ \langle \circ, \ell_1 \rangle \} \rangle : \\
&\quad \langle \ell_1, L_1 \rangle \in S_A[\text{sexp}_1](a, v), \langle \ell_2, L_2 \rangle \in S_A[\text{sexp}_2](a, v) \} \\
S_A[\text{sub}(sexp, k_1, k_2)](a, v) &= \{ \langle \ell_1, L_1 \cup \{ \langle \text{sub}, k_1, k_2 \rangle, \ell_1 \} \rangle : \langle \ell_1, L_1 \rangle \in S_A[\text{sexp}](a, v) \} \\
S_A[\text{hash}(sexp)](a, v) &= \{ \langle \ell_1, L_1 \cup \langle \text{hash}, \ell_1 \rangle \rangle : \langle \ell_1, L_1 \rangle \in S[\text{sexp}](a, v) \}
\end{aligned}$$

**Fig. 2.** Semantics of Expressions on Atomic Data

$$\begin{aligned}
S[x := sexp](a, v) &= (a[x \mapsto S_A[\text{sexp}](a, v)], v[x \mapsto S_S[\text{sexp}](v)]) \\
S[\text{send}(sexp)](a, v) &= (a, v) \\
S[c_1; c_2](a, v) &= S[c_2](S[c_1](a, v)) \\
S[\text{if } bexp \text{ then } c_1 \text{ else } c_2](a, v) &= \begin{cases} S[c_1](a, v) & \text{if } S_B[bexp](v) \\ S[c_2](a, v) & \text{otherwise} \end{cases} \\
S[\text{while } bexp \text{ do } c](a, v) &= S[\text{if } (bexp) \text{ (c; while } bexp \text{ do } c)](a, v)
\end{aligned}$$

**Fig. 3.** Concrete Semantics of Statements

Given a program  $p$ , we will denote the concrete datastore associated with this program by  $C_p$ .

*Example* Consider the `Inmobi` example from Section 2. Method `updateInfo` accesses various data coming from the datastore. We represent by (i)  $\langle \text{Language}, \emptyset \rangle$  the language returned by the `Default` object (line 30), (ii)  $\langle \text{Country}, \emptyset \rangle$  the country returned by the `Default` object (line 31), and (iii)  $\langle \text{AndroidId}, \emptyset \rangle$  the Android identifier. These three data sources are stable, that is, they always return the same values. For the locations (that is, line 34) it may be the case that different calls of `getLastKnownLocation` retrieves different locations. Therefore, the concrete datastore contains  $\langle \text{Location}_i, \emptyset \rangle : i \in \mathbb{N}$  as well. Instead, for the `IMSI` example we have only one datum  $\langle \text{IMSI}, \emptyset \rangle$ .

## 4.2 Semantics

We suppose that a standard concrete evaluation of numerical ( $S_N : nexp \times V \rightarrow \mathbb{Z}$ ) and string ( $S_S : sexp \times V \rightarrow \mathbb{S}$ ) expressions is provided, as well as the evaluation of Boolean conditions ( $S_B : bexp \times V \rightarrow \{\text{true}, \text{false}\}$ ). In addition, we suppose that the semantic evaluation of label expressions ( $S_L : lexp \times \Sigma \rightarrow \text{Lab}$ ) returns a data label given a label expression.

The evaluation of the expressions on atomic data  $S_A : sexp \times \Sigma \rightarrow \wp(\text{ID})$  is defined in Fig. 2. Observe that this enhanced concrete semantics of expressions can be seen as an abstract representation of partial execution traces, where each expression tree is projected to the data associated to the labels in `Lab`.

Once the semantics of expression is formalized, the (concrete enhanced) semantics of statements can be expressed as depicted in Figure 3.

*Example.* Consider again the `Inmobi` example of Section 2. After the execution of `updateInfo` (line 4) we have that (i) `user.language`  $\mapsto \{ \langle \text{Language}, \emptyset \rangle \}$ ,

(ii) `user.country`  $\mapsto \{\langle \text{Country}, \emptyset \rangle\}$ , (iii) `user.id`  $\mapsto \{\langle \text{AndroidId}, \{\langle \text{hash.Android-Id} \rangle\} \rangle\}$ , and (iv) `user.loc`  $\mapsto \{\langle \text{Location}_1, \emptyset \rangle\}$ . We then concatenate all this data in a string stored in `url` at line 17. Therefore, we obtain the following atomic data expression with label `AndroidId`:

$\langle \text{AndroidId}, \{(\text{hash}, \text{AndroidId}), (\circ, \text{Language}), (\circ, \text{Country}), (\circ, \text{Location}_1)\} \rangle$  while for `Location1` we obtain  $\langle \text{Location}_1, \{(\circ, \text{AndroidId})\} \rangle$  since this is the last element concatenated when building `url`. For the IMSI example, we obtain that the data expression leaked at line 2 is  $\langle \text{IMSI}, \{(\text{sub}, 0, 6), \text{IMSI}\} \rangle$ .

### 4.3 Canonical Form of Atomic Data

The definition of atomic data does not impose any constraint on the number of elements. In particular, the same source label can appear several times in an atomic datum, when its data have multiple impact on the expression's value. However, if we are just interested to observe the sources of an expression, and the set of operators applied to each source, a more compact representation of atomic data can be given, where each source label appears at most once.

Given an atomic datum  $d = \{\langle \ell_j, L_j \rangle : j \in J\}$ , we denote by  $\text{src}(d)$  its source set  $\{\ell_j : j \in J\}$ . Moreover, given a label  $\ell$  and an atomic datum  $d$ , we denote by  $\text{links}(\ell, d)$  the links set of  $\ell$  in  $d$  if  $\ell \in \text{src}(d)$ , and  $\emptyset$  otherwise.

**Definition 2.** *We say that an atomic datum  $d$  is in canonical form if every label in  $\text{src}(d)$  occurs as a source label exactly once in  $d$ . Given an atomic datum, its canonical form can be obtained by applying the following unary source collapse operator  $\rho$ :*

$$\rho(d) = \{\langle \ell, \cup \text{links}(\ell, d) \rangle : \ell \in \text{src}(d)\}$$

## 5 Abstract Semantics

There are two main ways to get abstractions of the concrete semantics defined so far: abstracting values, and abstracting labels. The abstract elements should be an overapproximation of the concrete values assigned to variables in the concrete computation steps.

### 5.1 Values Abstraction

Values can be abstracted by means of well-known either relational or non-relational domains for numerical and textual values [9,26]. Therefore, we suppose that a value abstract domain  $V^a$  is provided, and it is equipped with the standard lattice and semantic operators.

### 5.2 Labels Abstraction

Labels can be abstracted by any abstract domain for categorical data, like a flat constant propagation domain. Observe that when dealing with data stored in relational form, i.e. by means of bi-dimensional tables, a relational abstract domain for array representation can be adopted, as defined in [11].

*Example.* In the Inmobi example of Section 2, we do not need to apply any abstraction on Language, Country, and AndroidId, since these are persistent throughout the execution. Instead, we need to apply abstraction to the locations  $\text{Location}_i : i \in \mathbb{N}$ , since the statement at line 37 may produce many values. Therefore, we abstract together all the locations produced by the same program point pp with  $\text{Location}^{\text{pp}}$ . In our example, this means that we abstract the data source at line 37 with  $\text{Location}^{37}$ .

### 5.3 Atomic Data Abstraction

We are now in position to formalize the atomic data abstract domain  $\mathbb{AD}$ .

**Definition 3 (Abstract Atomic Data).** *Given a set of atomic data, an abstract element will be a set of tuples  $\{\langle \ell_j^a, L_j^{a\sqsupset}, L_j^{a\sqcup} \rangle : j \in J\} \in \mathbb{AD}$ , where*

- $\ell_j^a$  is an element of an abstract domain that abstracts labels in  $\text{Lab}$
- $L_j^{a\sqsupset} = \{\langle \text{op}_{ij}^a, \ell_{ij}^a \rangle : i \in I\}$  is an under-approximation of the set of operators applied to the sources represented by  $\ell_j^a$  with values coming from sources represented by  $\ell_{ij}^a$
- $L_j^{a\sqcup} = \{\langle \text{op}_{ij}^a, \ell_{ij}^a \rangle : i \in I'\}$  is an over-approximation of the set of operators applied to the sources represented by  $\ell_j^a$  with values coming from sources represented by  $\ell_{ij}^a$
- $L_j^{a\sqsupset} \subseteq L_j^{a\sqcup}$ .

The order on the abstract elements is given by the order on the Cartesian product of the components' domain, and the least upper bound and greatest lower bound operators are defined accordingly.

**Definition 4 (Partial Order on Abstract Atomic Data).** *Given two abstract atomic data  $d_1 = \{\langle \ell_{1i}^a, L_{1i}^{a\sqsupset}, L_{1i}^{a\sqcup} \rangle : i \in I_1\}$  and  $d_2 = \{\langle \ell_{2i}^a, L_{2i}^{a\sqsupset}, L_{2i}^{a\sqcup} \rangle : i \in I_2\}$  on the same abstract domains for values and labels,*

$$d_1 \sqsubseteq d_2 \Leftrightarrow \forall i \in I_1 \exists j \in I_2 : \ell_{1i}^a = \ell_{2j}^a, L_{1i}^{a\sqsupset} \supseteq L_{2j}^{a\sqsupset}, L_{1i}^{a\sqcup} \subseteq L_{2j}^{a\sqcup}$$

Given an abstract atomic datum  $\{\langle \ell_j^a, L_j^{a\sqsupset}, L_j^{a\sqcup} \rangle : j \in J\}$ , we denote by  $\text{src}(d)$  its source set  $\{\ell_j^a : j \in J\}$ .

**Definition 5 (Least Upper bound of Abstract Atomic Data).** *Given two abstract atomic data  $d_1 = \{\langle \ell_{1i}^a, L_{1i}^{a\sqsupset}, L_{1i}^{a\sqcup} \rangle : i \in I_1\}$  and  $d_2 = \{\langle \ell_{2i}^a, A, L_{2i}^{a\sqsupset}, L_{2i}^{a\sqcup} \rangle : i \in I_2\}$  on the same abstract domains for values and labels, the least upper bound of  $d_1$  and  $d_2$  is the atomic datum*

$$d_1 \sqcup d_2 = \bigcup_{\ell^a \in \text{src}(d_1) \cup \text{src}(d_2)} \left\{ \begin{array}{ll} \langle \ell^a, L_1^{a\sqsupset}, L_1^{a\sqcup} \rangle & \text{if } \ell^a \in \text{src}(d_1) \setminus \text{src}(d_2) \\ \langle \ell^a, L_2^{a\sqsupset}, L_2^{a\sqcup} \rangle & \text{if } \ell^a \in \text{src}(d_2) \setminus \text{src}(d_1) \\ \langle \ell^a, L_1^{a\sqsupset} \cap L_2^{a\sqsupset}, L_1^{a\sqcup} \cup L_2^{a\sqcup} \rangle & \text{otherwise} \end{array} \right\}$$

Let  $\text{Lab}^a$  and  $A$  be complete lattices featuring Galois Connections with the concrete domains of labels and values, respectively. Let  $(\gamma_{\text{Lab}}, \alpha_{\text{Lab}})$ ,  $(\gamma_A, \alpha_A)$  be the corresponding concretization and abstraction functions. When applied to a set of links  $\{\langle \text{op}_i, \ell_i \rangle : i \in I\}$ , the function  $\alpha_{\text{Lab}}$  returns the set  $\{\langle \text{op}_i^a, \alpha_{\text{Lab}}(\ell_i) \rangle : i \in I\}$ , where  $\text{op}_i^a$  is the abstract operator that safely approximates  $\text{op}_i$  in the abstract domain  $A$ .



**Definition 6 (Abstraction function).** *The abstraction function  $\alpha : \wp(\mathbb{D}) \longrightarrow \mathbb{AD}$  is first defined on singletons and then extended to sets by applying the least upper bound operator.*

$$\begin{aligned}\alpha_s(\{\langle \ell_i, L_i \rangle : i \in I\}) &= \{\langle \alpha_{\text{Lab}}(\ell_i), \alpha_{\text{Lab}}(L_i), \alpha_{\text{Lab}}(L_i) \rangle : i \in I\} \\ \alpha(\{d_j \in \mathbb{D} : j \in J\}) &= \bigsqcup_{j \in J} \alpha_s(d_j).\end{aligned}$$

Notice that in the definition above, when considering a single atomic datum, in its abstract representation the under- and over-approximations of the link sets are equal. The gap among these sets is introduced in fact by the least upper bound operator.

**Definition 7 (Concretization function).** *The concretization of abstract atomic data is defined as an adjoint of the abstraction function:  $\gamma(ad) = \{d \in \mathbb{D} : \alpha(d) \sqsubseteq ad\}$*

**Theorem 1 (Galois Connection).** *The functions  $\alpha$  and  $\gamma$  defined above form a Galois Connection between  $\wp(\mathbb{D})$  and  $\mathbb{AD}$ , i.e.:*

- i)  $\alpha$  and  $\gamma$  are monotone,
- ii)  $\forall ad \in \mathbb{AD} : \alpha(\gamma(ad)) \sqsubseteq_{\mathbb{AD}} ad$
- iii)  $\forall S \subseteq \mathbb{D} : S \subseteq \gamma(\alpha(S))$ .

*Proof.*  $\mathbb{AD}$  is the Cartesian product of abstract domains featuring Galois Connections with the concrete domain  $\wp(\mathbb{D})$ , and the functions  $\alpha$  and  $\gamma$  are defined in canonical way w.r.t. the Cartesian product [8].

We define by  $AD^a : \text{Var} \mapsto \wp(\mathbb{AD})$  the component of the abstract domain tracking information on atomic data expressions. The partial order, the upper bound and the concretization function are defined as pointwise application of the operators defined on  $\mathbb{AD}$ .

## 5.4 Abstract Domain

Our abstract domain is the Cartesian product of the Atomic Data abstract domain ( $AD^a$ ), and the value domain ( $V^a$ ).

## 5.5 Data-Store Abstraction

The analysis of a program  $P$  aimed at verifying that it satisfies a given confidentiality policy with respect to data stored in devices running  $P$  can be defined either as a "datastore-aware" analysis, i.e. running the analysis on the actual data contained in the device, or in a "datastore-unaware" way, i.e. running the analysis on a generic datastore that represents the actual datastores under a suitable abstraction of labels and values.

A "datastore-aware" analysis has the advantage of being in general more accurate, as it allows to deal with the actual values that are leaked by the program. However, this scenario requires the analysis being applied only once the program is installed on the device, as an app that runs on the device itself or on a third-party verifier that should be given access permission to the device's datastore.

The accuracy of a "datastore-unaware" analysis heavily relies on the datastore abstraction, but it has the advantage of being applicable to the program with no need to access the actual confidential data when running the analysis itself.

$$\begin{aligned}
S^a \llbracket x := \text{sexp} \rrbracket (a^a, v^a) &= (a^a, S_v^a \llbracket x := \text{sexp} \rrbracket (v^a)) \\
S^a \llbracket \text{send}(\text{sexp}) \rrbracket (a^a, v^a) &= (a^a, v^a) \\
S^a \llbracket c_1; c_2 \rrbracket (a^a, v^a) &= S^a \llbracket c_2 \rrbracket (S^a \llbracket c_1 \rrbracket (a^a, v^a)) \\
S^a \llbracket \text{if } \text{bexp} \text{ then } c_1 \text{ else } c_2 \rrbracket (a^a, v^a) &= S^a \llbracket c_1 \rrbracket (a^a, S_B^a \llbracket \text{bexp} \rrbracket (v^a)) \sqcup S^a \llbracket c_2 \rrbracket (a^a, S_B^a \llbracket \neg \text{bexp} \rrbracket (v^a)) \\
S^a \llbracket \text{while } \text{bexp} \text{ do } c \rrbracket (a^a, v^a) &= \text{fix}(S^a \llbracket \text{if } (\text{bexp}) \text{ while } \text{bexp} \text{ do } c \rrbracket (a^a, v^a))
\end{aligned}$$

**Fig. 4.** (Abstract) Semantics of Statements

**Definition 8 (Abstract Datastore).** Given a Galois connection between  $\wp(\mathbb{D})$  and  $\mathbb{AD}$ , an abstract datastore is a set  $D^a = \{\{\langle \ell_j^a, \emptyset, \emptyset \rangle\} : j \in J\} \subseteq \mathbb{AD}$  such that  $\forall i, j \in J : i \neq j \Rightarrow \ell_i^a \neq \ell_j^a$ .

An abstract datastore  $D^a = \{\{\langle \ell_j^a, \emptyset, \emptyset \rangle\} : j \in J\}$  is an abstraction of all concrete datastores  $D = \{\{\langle \ell_i, \emptyset \rangle\} : i \in I\}$  that satisfies the following conditions:  $\bigcup_{i \in I} \{\alpha(\ell_i)\} = \bigcup_{j \in J} \{\ell_j^a\}$ , and  $\forall i \in I \exists! j \in J$  such that  $\ell_i \in \gamma_{\text{Lab}}(\ell_j^a)$ .

*Example.* Consider the Inmobi example introduced in Section 2. In particular, we have the following atomic data expressions:  $\langle \text{Language}, \emptyset, \emptyset \rangle$ ,  $\langle \text{Country}, \emptyset, \emptyset \rangle$ , and  $\langle \text{AndroidId}, \emptyset, \emptyset \rangle$ . The only case that is slightly different regards Location since we abstract with  $\langle \text{Location}^{37}, \emptyset, \emptyset \rangle$  all the concrete data expressions in  $\langle \text{Location}_i, \emptyset \rangle : i \in \mathbb{N}$ . Similarly, for the IMSI example we have  $\langle \text{IMSI}, \emptyset, \emptyset \rangle$ .

## 5.6 Abstract Semantics of Statements

Figure 4 depicts the abstract semantics of statements<sup>5</sup>. We omit the abstract semantics of expressions, as it can be easily formalized by mimicking the concrete semantics, the only difference being that (1) every operation has impact on both link sets associated to an abstract label, and (2) abstract atomic data are kept in canonical form by systematically applying the following normalization operator in presence of multiple occurrences of the same label in the source set of an abstract atomic datum:  $\rho(ad) = \{\langle \ell^a, \cap \text{links}(\ell^a, ad), \cup \text{links}(\ell^a, ad) \rangle : \ell^a \in \text{src}(ad)\}$

*Example.* The abstract semantics of the Inmobi example does not substantially differ from the concrete semantics for the the example of Section 2. After the execution of `updateInfo` (line 4), we have the same information described in Section 4.2, the only difference being that the abstract label for the location is `Location`<sup>37</sup> instead of the concrete label `Location1`<sup>6</sup>. The same consideration applies to the results of the concatenation at line 17.

For the IMSI example, the abstract semantics tracks that we log the abstract datum  $\langle \text{IMSI}, \{([sub, 0, 6], \text{IMSI})\}, \{([sub, 0, 6], \text{IMSI})\} \rangle$ .

The following theorem formalizes the soundness of the analysis.

<sup>5</sup> Observe that this semantics does not capture indirect information flow.

<sup>6</sup> For the sake of simplicity, we ignore the issues related with heap abstraction. It has been demonstrated [16] that value domains (like  $AD^a$  and  $V^a$ ) can be combined with heap abstractions relying on standard operators of value domains.

**Theorem 2.** *The abstract semantics of a program  $P$  with an abstract datastore  $D^a$  is a conservative (sound) over-approximation of the enhanced concrete semantics of  $P$  with a concrete datastore  $D \in \gamma(D^a)$ .*

*Proof.* By induction on the length of the trace as in [23], by lifting the local correctness of the operations to the Cartesian product [8].

We observe in particular that for each execution of  $P$  with input  $\mathbf{I}$ , if a value assigned to a variable  $v$  in the store is obtained from values coming from local data stored in  $\ell$  through operations in  $R$ , then there is a corresponding abstract trace of  $P$  with input  $\alpha(\mathbf{I})$ , assigning  $v$  an abstract atomic datum  $ad$  such that (i)  $\ell^a \in \text{src}(ad)$ , and (ii)  $R \subseteq \{\text{op} : (\text{op}^a, \ell_j^a) \in \text{links}(\ell^a, ad)\}$  (where  $\ell^a = \alpha(\ell)$  is the label in the abstract datastore representing  $\ell$ ).

## 6 Confidentiality and Obfuscation

So far, we made no distinctions among data contained in the data-store, with respect to their confidentiality level. In general, we can consider a lattice of confidentiality levels  $\mathcal{S}$ , and we can associate to each label  $\ell$  in  $\text{Lab}$  an element  $s_\ell \in \mathcal{S}$ . Confidentiality levels are assigned to labels, and values corresponding to these labels will inherit from them the same confidentiality level.

On the operation side, we introduce the notion of *obfuscation degree*. The intuitive idea is that if you know which operation has been applied to get an expression, and the expression itself, you can look at the amount of information which is necessary to recover the sources the operation applied to. This leads us to assume the existence of a partial-order relation among operations that captures their different obfuscation impact.

This can be seen as a generalization of the all-or-nothing tainting approach [37,39], where only declassification operators (e.g., encryption) are tracked.

The obfuscation degree of an operator can be seen as a measure of the complexity of the brute-force analysis needed by an external observer in order to detect the actual source data when knowing just the result of the operation and the applied operator.

**Definition 9 (Obfuscation Degree).** *Consider a complete lattice  $(\mathbb{O}, \sqsubseteq_{\mathbb{O}})$ , and a map  $\zeta : \text{Op} \rightarrow \mathbb{O}$ , such that  $\zeta(\text{op}_1) \sqsubseteq_{\mathbb{O}} \zeta(\text{op}_2)$  if the obfuscation power of  $\text{op}_1$  is smaller than the obfuscation power of  $\text{op}_2$ . We say that the obfuscation degree of an operator  $\text{op} \in \text{Op}$  is  $\zeta(\text{op})$ .*

*Example.* The string operators of *sexp* introduced in Section 3 have different obfuscation degrees. For instance, *encrypt* obfuscates more than *hash*, while the power of obfuscation of *substring* may depend on the indexes used to compute the substring, and the particular information contained in the string. For instance, in the IMSI example of Section 2 the substring operator at line 2 has a high obfuscation degree, but this relies on the value information tracked on the indexes passed to *substring*.

### 6.1 Confidentiality of Atomic Data

Given an atomic datum, a confidentiality value can be assigned to it by considering an under- and over-approximation of the confidentiality levels of source data, and by

considering an under- and over- approximation of obfuscation power of the operations applied to them.

We first define it at a concrete level, on top of our instrumented atomic data semantics, and then we can lift this notion to the abstract case.

**Definition 10 (Confidentiality of Atomic Data for Monotonic Operators).** *Let  $S$  be a lattice representing confidentiality levels of labels. Let  $\mathbb{O}$  be a lattice representing the obfuscation power of operators. Finally, let  $\eta$  and  $\zeta$  be functions assigning confidentiality/obfuscation values in  $S$  and  $\mathbb{O}$  to labels and operators, respectively.*

*If the combination of operators in  $\bigcup_{i \in I} L_i$  is monotonic with respect to the obfuscation order in the lattice  $\mathbb{O}$ , the confidentiality value of an atomic datum  $\{\langle \ell_i, L_i \rangle : i \in I\}$  with respect to  $(\eta, \zeta)$  is the tuple  $(sc_{min}, sc_{max}, lc_{min}, lc_{max})$ , where:*

$$\begin{aligned} sc_{min} &= \sqcap_{\mathbb{D}} \{\eta(\ell_i) : i \in I\} & lc_{min} &= \sqcap_{\mathbb{O}} \{\zeta(\{op_{ij} : (op_{ij}, \ell_j) \in L_i, i \in I\})\} \\ sc_{max} &= \sqcup_{\mathbb{D}} \{\eta(\ell_i) : i \in I\} & lc_{max} &= \sqcup_{\mathbb{O}} \{\zeta(\{op_{ij} : (op_{ij}, \ell_j) \in L_i, i \in I\})\} \end{aligned}$$

*Example.* Imagine that we have  $L < M < H$  as both the confidentiality and obfuscation lattice. We then establish that *encrypt* has  $H$  obfuscation level, and *hash* has level  $M$ , whereas the obfuscation level of substring depends on the parameters:  $[sub, k_1, k_2]$  has level  $L$  if  $k_1 = 6$  and  $k_2 = 9$ , it has level  $M$  if  $6 < k_1 + k_2 < 15$ , and it has level  $H$  if  $k_1 + k_2 \leq 6$ . Consider then the concrete labels of the *Inmobi* example introduced in Section 4.1. We define as  $L$  both *Language* and *Country*, since they do not contain particularly confidential information. Instead, we define as  $H$  *AndroidId*, since this datum allows to uniquely identify our Android account, and track our activity. Finally,  $Location_i : i \in \mathbb{N}$  are all  $M$ , since these locations allow to identify our geographical location at a given point, but do not uniquely identify us. For the *IMSI* example, we obtain that for the data expression leaked at line 2, i.e.  $\langle IMSI, \{([sub, 0, 6], IMSI)\}$ , we get  $sc_{min} = sc_{max} = H$  and  $lc_{min} = lc_{max} = H$ . This says that even if sensitive data items are leaked, a powerful obfuscation is definitely applied to them before releasing them.

Notice that Definition 10 is explicitly restricted to the case of operators whose combination is monotonic with respect to the obfuscation order in  $\mathbb{O}$ . If we are interested to consider also programs where the combination of operators is non-monotonic, we just need to give an obfuscation value to sets of operators instead of singletons.

**Definition 11 (Confidentiality of Atomic Data - General Case).** *Let  $S$  be a lattice representing confidentiality levels of labels. Let  $\mathbb{O}$  be a lattice representing the obfuscation power of operators. Finally, let  $\eta$  be a function assigning confidentiality values in  $S$  to labels, and let  $\zeta$  be a function assigning to each set of operators an interval in  $\mathbb{O} \times \mathbb{O}$  representing its min and max obfuscation power.*

*The confidentiality value of an atomic datum  $\{\langle \ell_i, L_i \rangle : i \in I\}$  with respect to  $(\eta, \zeta)$  is the tuple  $(sc_{min}, sc_{max}, lc_{min}, lc_{max})$ , where:*

$$\begin{aligned} sc_{min} &= \sqcap_{\mathbb{D}} \{\eta(\ell_i) : i \in I\} & lc_{min} &= \sqcap_{\mathbb{O}} \{\pi_1(\zeta(\{op_{ij} : (op_{ij}, \ell_j) \in L_i\})) : i \in I\} \\ sc_{max} &= \sqcup_{\mathbb{D}} \{\eta(\ell_i) : i \in I\} & lc_{max} &= \sqcup_{\mathbb{O}} \{\pi_2(\zeta(\{op_{ij} : (op_{ij}, \ell_j) \in L_i\})) : i \in I\} \end{aligned}$$

*where  $\pi_1$  and  $\pi_2$  denote the on the min and max element of the interval, respectively.*

Notice that keeping track of minimal confidentiality and maximal obfuscation allows us (when they are equal to maximal confidentiality and minimal obfuscation, respectively) to be aware of the precision of these values.

## 6.2 Confidentiality of Abstract Atomic Data

In order to define the confidentiality value of abstract atomic data, we need to assign a confidentiality value to abstract labels. As abstract labels may represent concrete labels with different confidentiality values, the confidentiality function  $\eta^a$  returns an interval min/max of values in  $\mathbb{S} \times \mathbb{S}$  instead of a single value. On the obfuscation side, we just lift the value, as we can assume there is always a one-to-one correspondence between concrete and abstract operators.

As in the concrete setting, we distinguish the case in which all operator combinations behave monotonically with respect to the obfuscation order, from the general case, that takes into account non-monotonic behaviors, the price to pay being to assign obfuscation values to sets of operators instead of single ones.

**Definition 12 (Confidentiality Value of Abstract Atomic Data - Monotonic Operators).** *Let  $\mathbb{S}$  and  $\mathbb{O}$  be the lattices representing the labels' confidentiality and the obfuscation power of operators, respectively. Let  $\eta$  and  $\zeta$  be functions assigning confidentiality/obfuscation values in  $\mathbb{S}$  and  $\mathbb{O}$  to (concrete) labels and operators, respectively. Let  $\eta^a : \text{Lab}^a \rightarrow \mathbb{S} \times \mathbb{S}$  such that  $\eta^a(\ell^a) = [\sqcap\{\eta(\ell) : \ell \in \gamma(\ell^a)\}, \sqcup\{\eta(\ell) : \ell \in \gamma(\ell^a)\}]$ . Let  $\zeta^a$  be the function assigning to each abstract operator the same obfuscation value assigned by  $\zeta$  to the concrete operator it corresponds.*

*If the combination of operators in  $\mathbb{O}_p$  appearing in  $\bigcup_{i \in I} L_i^{a\sqcup}$  is monotonic with respect to the obfuscation order in  $\mathbb{O}$ , then the confidentiality value of an abstract atomic datum  $\{(\ell_i^a, \text{val}_i^a, A_i, L_i^{a\sqcap}, L_i^{a\sqcup}) : i \in I\}$  is a tuple  $(sc_{min}^a, sc_{max}^a, lc_{min}^a, lc_{max}^a)$ , where:*

$$\begin{aligned} sc_{min}^a &= \sqcap_{\mathbb{S}}\{\pi_1(\eta^a(\ell_i^a)) : i \in I\} & lc_{min}^a &= \sqcap_{\mathbb{O}}\{\zeta^a(op_{ij}^a) : (op_{ij}^a, \ell_j^a) \in L_i^{a\sqcap}, i \in I\} \\ sc_{max}^a &= \sqcup_{\mathbb{S}}\{\pi_2(\eta^a(\ell_i^a)) : i \in I\} & lc_{max}^a &= \sqcup_{\mathbb{O}}\{\zeta^a(op_{ij}^a) : (op_{ij}^a, \ell_j^a) \in L_i^{a\sqcup}, i \in I\} \end{aligned}$$

*where  $\pi_1$  and  $\pi_2$  denote the min and max element of the interval, respectively.*

Observe that  $lc_{min}^a$  is obtained as the greatest lower bound of the obfuscation values corresponding to operators that are surely applied to compute the value, while  $lc_{max}^a$  is obtained as the least upper bound of the obfuscation values corresponding to operators that are possibly applied to compute the value. As  $L_i^{a\sqcap} \subseteq L_i^{a\sqcup}$  for each  $i \in I$ , we get that  $lc_{min}^a \sqsubseteq_{\mathbb{O}} lc_{max}^a$ .

**Definition 13 (Confidentiality Value of Abstract Atomic Data - General Case).** *Let  $\mathbb{S}$  and  $\mathbb{O}$  be the lattices representing the labels' confidentiality and the obfuscation power of operators, respectively. Let  $\eta$  and  $\zeta$  be functions assigning confidentiality/obfuscation values in  $\mathbb{S}$  and  $\mathbb{O}$  to (concrete) labels and operators, respectively. Let  $\eta^a : \text{Lab}^a \rightarrow \mathbb{S} \times \mathbb{S}$  such that  $\eta^a(\ell^a) = [\sqcap\{\eta(\ell) : \ell \in \gamma(\ell^a)\}, \sqcup\{\eta(\ell) : \ell \in \gamma(\ell^a)\}]$ . Finally, let  $\zeta^a$  be a function assigning to each set of (abstract) operators an interval in  $\mathbb{O} \times \mathbb{O}$  representing the min and max obfuscation power.*

*The confidentiality value of an abstract atomic datum  $\{(\ell_i^a, \text{val}_i^a, A_i, L_i^{a\sqcap}, L_i^{a\sqcup}) : i \in I\}$  is a tuple  $(sc_{min}^a, sc_{max}^a, lc_{min}^a, lc_{max}^a)$ , where:*

$$\begin{aligned} sc_{min}^a &= \sqcap_{\mathbb{S}}\{\pi_1(\eta^a(\ell_i^a)) : i \in I\} & lc_{min}^a &= \sqcap_{\mathbb{O}}\{\pi_1(\zeta^a(S)) : S \subseteq \{op_{ij}^a : (op_{ij}^a, \ell_j^a) \in L_i^{a\sqcup}\}, i \in I\} \\ sc_{max}^a &= \sqcup_{\mathbb{S}}\{\pi_2(\eta^a(\ell_i^a)) : i \in I\} & lc_{max}^a &= \sqcup_{\mathbb{O}}\{\pi_2(\zeta^a(S)) : S \subseteq \{op_{ij}^a : (op_{ij}^a, \ell_j^a) \in L_i^{a\sqcup}\}, i \in I\} \end{aligned}$$

*where  $\pi_1$  and  $\pi_2$  denote the on the min and max element of the interval, respectively.*

Notice that  $lc_{min}^a$  and  $lc_{max}^a$  in the general case are both obtained as the greatest lower bound and least upper bound, respectively, of the obfuscation values corresponding to operators that are possibly applied to compute the value. This conservative approach guarantees the soundness of the result also in presence of operators whose combination does not behave monotonically with respect to obfuscation, i.e. the monotonicity of confidentiality with respect to the partial order in the domain of abstract atomic data.

## 7 Privacy Compliance Policies

**Definition 14 (Confidentiality Policy).** *Given the set of data source labels  $\text{Lab}$ , and the confidentiality/obfuscation lattices  $\mathbb{S}$  and  $\mathbb{O}$  for labels and operations, respectively, a confidentiality policy is a tuple  $\pi = (\eta, \zeta, \kappa_{sc_{max}}, \kappa_{lc_{min}})$  such that*

- $\eta, \zeta$  assign each label and each operator a corresponding value in the confidentiality lattices  $\mathbb{S}$  and  $\mathbb{O}$ , respectively.
- $\kappa_{sc_{max}}$  is a source confidentiality threshold (the max confidentiality level allowed for sources).
- $\kappa_{lc_{min}}$  is an obfuscation threshold (the min obfuscation level required for operators).

Given a program  $P$ , let  $X$  be the set of concrete/abstract atomic data  $P$  generated as an output. We say that  $P$  satisfies the confidentiality policy  $\pi = (\eta, \zeta, \kappa_{sc_{max}}, \kappa_{lc_{min}})$  if:

$\forall d \in X$ , if  $(sc_{min}, sc_{max}, lc_{min}, lc_{max})$  is the confidentiality value of  $d$  with respect to  $(\eta, \zeta)$ , then,  $sc_{max} \sqsubseteq_{\mathbb{S}} \kappa_{sc_{max}}$  and  $lc_{min} \sqsupseteq_{\mathbb{O}} \kappa_{lc_{min}}$ .

**Theorem 3.** *Consider a program  $P$ , an abstract datastore  $A$ , and a confidentiality policy  $\pi = (\eta, \zeta, \kappa_{sc_{max}}, \kappa_{lc_{min}})$ . If the program  $P$  terminates, and the output of the analysis on  $P$  and  $A$  satisfies  $\pi$ , then any actual execution of program  $P$  on a concrete datastore in  $\gamma(A)$  satisfies the confidentiality policy  $\pi$ .*

*Proof.* By Theorem 2, and by the monotonicity of confidentiality values with respect to the partial order on the domain of abstract atomic data.

*Example.* A reasonable privacy policy for the `Inmob1` example of Section 2 may be that a datum can be released only if its obfuscation level is equal or higher than its confidentiality level. This program satisfies this model for `Country` and `Language` (whose confidentiality level is `L` and they are released without any obfuscation), but not for `Location1` (with confidentiality level `M` and released without any obfuscation) and `AndroidId` (whose confidentiality level is `H` and it is released after invoking `hash`, that is, with obfuscation level `M`).

### 7.1 Sources' Confidentiality Policies

The definition of confidentiality policy of atomic data discussed in the previous sections allows to capture the min and max levels of confidentiality/obfuscation carried by the values returned by a program, or shared with other applications.

As an orthogonal approach, we may define (at the concrete level) a confidentiality policy as a propositional formula that captures constraints on the allowed releasing levels of confidential data in the datastore, and verify if the atomic data returned by the concrete execution of the program satisfy that formula.

Let  $\text{Lab}$  denote as usual the set of data source labels, and  $\text{Op}$  denote the set of operators in the program. Consider a set of propositional variables  $V$  (with empty intersection with the set of program variables), and a function  $\lambda$  that maps elements of  $V$  into either labels or links.

$$\lambda : V \longrightarrow \text{Lab} \cup \{(op, \ell) : op \in \text{Op}, \ell \in \text{Lab}\}$$

A *policy formula* is a positive propositional formula on  $V$ , i.e. a propositional formula using only  $\wedge, \vee$  and  $\leftrightarrow$  logical operators<sup>7</sup>.

*Example.* For instance, we can express the fact that we can leak the Android ID if encrypted and the location, or the hashed Android ID, or the first six characters of the IMSI, by means of the formula  $\varphi = (x \wedge y) \vee z \vee w$ , where  $\lambda(x) = ([\text{encrypt}, k], \text{AndroidId})$ ,  $\lambda(y) = \text{Location}$ , and  $\lambda(z) = (\text{hash}, \text{AndroidId})$ ,  $\lambda(w) = ([\text{sub}, 0, 6], \text{AndroidId})$ .

Given a set of atomic data  $S$ , a set of propositional variables  $V$  and an assignment  $\lambda$  on  $V$ , we say that  $S$  satisfies the policy formula  $\varphi$  on  $V$  if  $S, \lambda \models \varphi$ , as defined inductively as follows:

$$\begin{aligned} S, \lambda \models v \in V &\Leftrightarrow \begin{cases} \lambda(v) \in \bigcup_{d \in S} \text{src}(d) & \text{if } \lambda(v) \in \text{Lab} \\ \lambda(v) \in \bigcup_{d \in S} \text{links}(d) & \text{otherwise.} \end{cases} \\ S, \lambda \models \varphi_1 \vee \varphi_2 &\Leftrightarrow S, \lambda \models \varphi_1 \text{ or } S, \lambda \models \varphi_2 \\ S, \lambda \models \varphi_1 \wedge \varphi_2 &\Leftrightarrow S, \lambda \models \varphi_1 \text{ and } S, \lambda \models \varphi_2 \\ S, \lambda \models \varphi_1 \leftrightarrow \varphi_2 &\Leftrightarrow S, \lambda \models \varphi_1 \text{ iff } S, \lambda \models \varphi_2. \end{aligned}$$

Observe that, by construction, if the data resource denoted by  $\lambda(x)$  contributes to any of the values represented by the atomic data  $S$ , then  $S, \lambda \models x$ .

*Example* The formula  $\varphi$  in the Example above is satisfied by the data expressions leaked by either `Inmob` and `IMSI` as described in the Example in Section 3.2.

Observe that checking the policy reduces to checking the satisfiability of the propositional assignment [42].

When we lift to the abstract setting, three-valued models of propositional formulas should be used in order to preserve soundness. Let  $S^a$  be a set of abstract atomic data, and  $\lambda^a : V \longrightarrow \text{Lab}^a \cup \{(op^a, \ell^a) : op \in \text{Op}, \ell^a \in \text{Lab}^a\}$  be a function mapping propositional variables into abstract labels and links. Consider the three value assignment

<sup>7</sup> The advantage of using positive formulas, i.e. formulas that are satisfied by assigning true to all its propositional variables, is that they well capture monotonic behaviors [42].

$assign(S^a, \lambda) : \mathbb{V} \rightarrow \{\mathbf{true}, \mathbf{false}, \top\}$  defined by

$$assign(S^a, \lambda^a)(v) = \begin{cases} \mathbf{true} & \text{if } \lambda^a(v) \in \mathbf{Lab}^a \text{ and } \lambda^a(v) \in \bigcup_{d \in S} \mathbf{links\_lab}(L_d^\square) \\ & \text{or} \\ & \text{if } \lambda^a(v) \notin \mathbf{Lab}^a \text{ and } \lambda^a(v) \in \bigcup_{d \in S} L^\square(d) \\ \mathbf{false} & \text{if } \lambda^a(v) \in \mathbf{Lab}^a \text{ and } \lambda^a(v) \notin \bigcup_{d \in S} \mathbf{src}(d) \\ & \text{or} \\ & \text{if } \lambda^a(v) \notin \mathbf{Lab}^a \text{ and } \lambda^a(v) \notin \bigcup_{d \in S} L^\square(d) \\ \top & \text{otherwise} \end{cases}$$

(where  $L^\square(d)$  and  $L^\square(d)$  denote the last two components of  $d$ , respectively), and consider the logical operators extended to the top value  $\top$  by:

$$\frac{\wedge \mid \mathbf{true} \ \mathbf{false} \ \top}{\top \mid \top \ \ \ \ \ \mathbf{false} \ \top} \quad \frac{\vee \mid \mathbf{true} \ \mathbf{false} \ \top}{\top \mid \mathbf{true} \ \top \ \ \ \ \ \top} \quad \frac{\leftrightarrow \mid \mathbf{true} \ \mathbf{false} \ \top}{\top \mid \top \ \ \ \ \ \top \ \ \ \ \ \top}$$

**Theorem 4.** *Let  $\varphi$  be a positive formula on a set  $\mathbb{V}$  of propositional variables,  $S^a$  be a set of abstract atomic data, and  $\lambda^a : \mathbb{V} \rightarrow \mathbf{Lab}^a \cup \{(op^a, \ell^a) : op \in \mathbf{Op}, \ell^a \in \mathbf{Lab}^a\}$  be a function mapping propositional variables into abstract labels and links. If  $assign(S, \lambda)(\varphi) = \mathbf{true}$ , then there is a set of atomic data  $S \subseteq \bigcup_{d \in S^a} \gamma(d)$  and a function  $\lambda$  satisfying  $\forall v \in \mathbb{V} : \lambda(v) \in \gamma_{\mathbf{Lab}}(\lambda^a(v))$ , such that  $S, \lambda \models \varphi$ .*

*Proof.* By structural induction on the formula  $\varphi$ , and by Theorem 1.

## 8 Related Work

In this section, we discuss in details how our work compares to some similar approaches in the area in addition to the high-level overview on the state of the art of Section 1.1.

Quantitative Information Flow (QIF) [25] is aimed at measuring the quantity of information that is leaked by a program. A given confidential datum might be manipulated by the program, that at the end releases only partial information. Then, the analysis checks if the quantity of released information is below a given threshold. Our approach shares with QIF the intuition that is crucial to estimate the quantity of information revealed, since often it is necessary to partially disclose a part of the information. Nevertheless, instead of measuring a quantity, our approach tracks the set of operators that have been applied to the datum before its release, and then we check if this matches what specified on the policy. We believe that QIF can be seen as an abstraction of a concrete semantics tracking the exact order of operators applied to a datum (instead of a set of operators as we do). In addition, QIF can track implicit flows, while we explicitly ignored these flows as we believe they lead to many false alarms.

Declassification-based approaches [33] suppose that a list of declassifier operators is given, and as soon as one of these operators is applied to a datum, then it can be sent to a sink. Our analysis can represent declassifiers through a policy stating that it is allowed to release the data on which at least one of these operators has been applied.



Decentralized Information Flow Control [20,27] systems represent a finer grained and more expressive model, in which each process can declassify information, rather than a central authority as in centralized system and classical declassification-based approaches. Nevertheless, our analysis can support this more complex scenarios by defining specific policies per process.

Another approach that has gained relevant results recently is Differential Privacy [13,14]. Given a data store, the goal of Differential Privacy is to discover if the variation of a query over the data set stays below a given threshold when an entry is added. Usually, some statistical noise is added in order to ensure differential privacy. Our approach may track how the data from the data set is aggregated and noise added, and this may be a first step towards proving differential privacy. However, how to relate this information on the operators applied to data and  $\epsilon$ -differential privacy is not straightforward at all, and it requires further investigation.

## 9 Conclusion

Our semantic framework for fine-grained information-flow analysis captures how the values released by an application may partially reveal confidential data stored on the device through different levels of obfuscation. The enhanced concrete semantics and the generic abstract domain we presented provide a workbench for (static) analysis of mobile apps that can be tuned by setting a few parameters: the domain representing values, the domain representing data locations, and the confidentiality and obfuscation values for data and operators. This data-centric approach may be utilized to refine existing tools like [15,41,40] aimed at enforcing privacy policies, providing the user with more accurate privacy control.

The problem of formalizing how the semantics of operations reflects on the corresponding obfuscation values, as well as the problem of assisting the user in the definition of privacy compliance policies remain of course, as there is a tradeoff between the amount of sensitive information that she allows the device to release and the accuracy and efficiency of some functionalities.

**Acknowledgments.** Work partially supported by PRIN "Security Horizons".

## References

1. Alvim, M.S., Scedrov, A., Schneider, F.B.: When not all bits are equal: Worth-based information flow. In: Abadi, M., Kremer, S. (eds.) POST 2014. LNCS, vol. 8414, pp. 120–139. Springer, Heidelberg (2014)
2. Amtoft, T., Banerjee, A.: A logic for information flow analysis with an application to forward slicing of simple imperative programs. *Science of Computer Programming* 64, 3–28 (2007)
3. AppBrain. Adnetwork stats, <http://www.appbrain.com/stats/libraries/ad>
4. Arzt, S., Rasthofer, S., et al.: Flowdroid: Precise context, flow, field, object-sensitive and lifecycle-aware taint analysis for android apps. In: PLDI. ACM (2014)
5. Askarov, A., Myers, A.: A semantic framework for declassification and endorsement. In: Gordon, A.D. (ed.) ESOP 2010. LNCS, vol. 6012, pp. 64–84. Springer, Heidelberg (2010)
6. Cavadini, S.: Secure slices of insecure programs. In: ASIACCS. ACM Press (2008)
7. Chaudhuri, A.: Language-based security on android. In: PLAS. ACM (2009)

8. Cortesi, A., Costantini, G., Ferrara, P.: A survey on product operators in abstract interpretation. *EPTCS* 129, 325–336 (2013)
9. Costantini, G., Ferrara, P., Cortesi, A.: Static analysis of string values. In: Qin, S., Qiu, Z. (eds.) *ICFEM 2011*. LNCS, vol. 6991, pp. 505–521. Springer, Heidelberg (2011)
10. Cousot, P., Cousot, R.: Abstract interpretation: Past, present and future. In: *CSL-LICS*. ACM (2014)
11. Cousot, P., Cousot, R., Logozzo, F.: A parametric segmentation functor for fully automatic and scalable array content analysis. In: *POPL*. ACM (2011)
12. Denning, D.E.: A lattice model of secure information flow. *Communications of the ACM* 19, 236–243 (1976)
13. Dwork, C.: Differential privacy: A survey of results. In: Agrawal, M., Du, D.-Z., Duan, Z., Li, A. (eds.) *TAMC 2008*. LNCS, vol. 4978, pp. 1–19. Springer, Heidelberg (2008)
14. Ebadi, H., Sands, D., Schneider, G.: Differential privacy: Now it's getting personal. In: *POPL*. ACM (2015)
15. Enck, W., Gilbert, P., et al.: Taintdroid: An information flow tracking system for real-time privacy monitoring on smartphones. *Comm. of the ACM* 57(3), 99–106 (2014)
16. Ferrara, P.: Generic combination of heap and value analyses in abstract interpretation. In: McMillan, K.L., Rival, X. (eds.) *VMCAI 2014*. LNCS, vol. 8318, pp. 302–321. Springer, Heidelberg (2014)
17. Halder, R., Zanioli, M., Cortesi, A.: Information leakage analysis of database query languages. In: *SAC*. ACM (2014)
18. Hammer, C., Snelting, G.: Flow-sensitive, context-sensitive, and object-sensitive information flow control based on program dependence graphs. *International Journal of Information Security* 8, 399–422 (2009)
19. Hornyack, P., Han, S., Jung, J., Schechter, S.E., Wetherall, D.: These aren't the droids you're looking for: retrofitting android to protect data from imperious applications. In: *CCS*. ACM (2011)
20. Krohn, M.N., Tromer, E.: Noninterference for a practical DIFC-based operating system. In: *IEEE S&P*. IEEE (2009)
21. Li, B.: Analyzing information-flow in java program based on slicing technique. *SIGSOFT Softw. Eng. Notes* 27, 98–103 (2002)
22. Lochbihler, A., Snelting, G.: On temporal path conditions in dependence graphs. *Journal of Automated Software Engineering* 16, 263–290 (2009)
23. Logozzo, F.: Class invariants as abstract interpretation of trace semantics. *Computer Languages, Systems & Structures* 35, 100–142 (2009)
24. Lu, L., Li, Z., Wu, Z., Lee, W., Jiang, G.: Chex: statically vetting android apps for component hijacking vulnerabilities. In: *CCS*. ACM (2012)
25. McCamant, S., Ernst, M.D.: Quantitative information flow as network flow capacity. In: *PLDI*. ACM (2008)
26. Miné, A.: The octagon abstract domain. *Higher-Order and Symbolic Computation* 19(1), 31–100 (2006)
27. Myers, A.C., Liskov, B.: A decentralized model for information flow control. In: *SOSP*. ACM (1997)
28. Nanevski, A., Banerjee, A., Garg, D.: Dependent type theory for verification of information flow and access control policies. *ACM TOPLAS* 35(2), 6:1–6:41 (2013)
29. Omoronyia, I., Cavallaro, L., et al.: Engineering adaptive privacy: on the role of privacy awareness requirements. In: *ICSE*. IEEE/ACM (2013)
30. Pottier, F., Simonet, V.: Information flow inference for ml. *ACM Transactions on Programming Languages and Systems* 25, 117–158 (2003)
31. Rasthofer, S., Lovat, E., Bodden, E.: Droid force: Enforcing complex, data-centric, system-wide policies in android. In: *ARES* (2014)

32. Sabelfeld, A., Myers, A.C.: Language-based information-flow security. *IEEE Journal on Selected Areas in Communications* 21, 5–19 (2003)
33. Sabelfeld, A., Sands, D.: Declassification: Dimensions and principles. *Journal of Computer Security* 17, 517–548 (2009)
34. Smith, G.: Principles of secure information flow analysis. In: Christodorescu, M., et al. (eds.) *Malware Detection. Advances in Information Security*, vol. 27, pp. 291–307. Springer (2007)
35. Smith, S.F., Thober, M.: Refactoring programs to secure information flows. In: *PLAS*. ACM (2006)
36. Spiekermann, S., Cranor, L.F.: Engineering privacy. *IEEE Trans. Software Eng.* 35(1), 67–82 (2009)
37. Sridharan, M., Artzi, S., Pistoia, M., Guarnieri, S., Tripp, O., Berg, R.: F4f: Taint analysis of framework-based web applications. In: *OOPSLA*. ACM (2011)
38. Tripp, O., Ferrara, P., Pistoia, M.: Hybrid security analysis of web javascript code via dynamic partial evaluation. In: *ISSTA*. ACM (2014)
39. Tripp, O., Pistoia, M., Fink, S.J., Sridharan, M., Weisman, O.: Taj: Effective taint analysis of web applications. In: *PLDI*. ACM (2009)
40. Tripp, O., Rubin, J.: A bayesian approach to privacy enforcement in smartphones. In: *USENIX Security* (2014)
41. Xiao, X., Tillmann, N., Fähndrich, M., de Halleux, J., Moskal, M.: User-aware privacy control via extended static-information-flow analysis. In: *ASE*. ACM (2012)
42. Zanioli, M., Ferrara, P., Cortesi, A.: Sails: Static analysis of information leakage with sample. In: *SAC*. ACM (2012)