

A Bayesian Approach to Privacy Enforcement in Smartphones

Omer Tripp
IBM Research, USA

Julia Rubin
IBM Research, Israel

Abstract

Mobile apps often require access to private data, such as the device ID or location. At the same time, popular platforms like Android and iOS have limited support for user privacy. This frequently leads to unauthorized disclosure of private information by mobile apps, e.g. for advertising and analytics purposes. This paper addresses the problem of privacy enforcement in mobile systems, which we formulate as a classification problem: When arriving at a privacy sink (e.g., database update or outgoing web message), the runtime system must classify the sink’s behavior as either legitimate or illegitimate. The traditional approach of information-flow (or taint) tracking applies “binary” classification, whereby information release is legitimate iff there is no data flow from a privacy source to sink arguments. While this is a useful heuristic, it also leads to false alarms.

We propose to address privacy enforcement as a learning problem, relaxing binary judgments into a quantitative/probabilistic mode of reasoning. Specifically, we propose a Bayesian notion of statistical classification, which conditions the judgment whether a release point is legitimate on the evidence arising at that point. In our concrete approach, implemented as the BAYESDROID system which is soon to be featured in a commercial product, the evidence refers to the similarity between the data values about to be released and the private data stored on the device. Compared to TaintDroid, a state-of-the-art taint-based tool for privacy enforcement, BAYESDROID is substantially more accurate. Applied to 54 top-popular Google Play apps, BAYESDROID is able to detect 27% [64] privacy violations with only 1 false alarm.

1 Introduction

Mobile apps frequently demand access to private information. This includes unique device and user identifiers,

such as the phone number or IMEI number (identifying the physical device); social and contacts data; the user’s location; audio (microphone) and video (camera) data; etc. While private information often serves the core functionality of an app, it may also serve other purposes, such as advertising, analytics or cross-application profiling [11]. From the outside, the user is typically unable to distinguish legitimate usage of their private information from illegitimate scenarios, such as sending of the IMEI number to a remote advertising website to create a persistent profile of the user.

Existing platforms provide limited protection against privacy threats. Both the Android and the iOS platforms mediate access to private information via a permission model. Each permission is mapped to a designated resource, and holds per all application behaviors and resource accesses. In Android, permissions are given or denied at installation time. In iOS, permissions are granted or revoked upon first access to the respective resource. Hence, both platforms cannot disambiguate legitimate from illegitimate usage of a resource once an app is granted the corresponding permission [10].

Threat Model In this paper, we address privacy threats due to authentic (as opposed to malicious) mobile applications [4, 19]. Contrary to malware, such applications execute their declared functionality, though they may still expose the user to unnecessary threats by incorporating extraneous behaviors — neither required by their core business logic nor approved by the user [13] — such as analytics, advertising, cross-application profiling, social computing, etc. We consider unauthorized release of private information that (almost) unambiguously identifies the user as a privacy threat. Henceforth, we dub such threats *illegitimate*.

While in general there is no bullet-proof solution for privacy enforcement, which can deal with any type of covert channel, implicit flow or proprietary data transformation, and even conservative enforcement approaches

can easily be bypassed [20], there is strong evidence that authentic apps rarely exhibit these challenges. According to a recent study [11], and also our empirical data (presented in Section 5), private information is normally sent to independent third-party servers. Consequently, data items are released in clear form, or at most following well-known encoding/encryption transformations (like Base64 or MD5), to meet the requirement of a standard and general client/server interface.

The challenge, in this setting, is to determine whether the app has taken sufficient means to protect user privacy. Release of private information, even without user authorization, is still legitimate if only a small amount of information has been released. As an example, if an application obtains the full location of the user, but then releases to an analytics server only coarse information like the country or continent, then in most cases this would be perceived as legitimate.

Privacy Enforcement via Taint Analysis The shortcomings of mobile platforms in ensuring user privacy have led to a surge of research on realtime privacy monitoring. The foundational technique grounding this research is *information-flow tracking*, often in the form of *taint analysis* [24, 17]: Private data, obtained via privacy *sources* (e.g. `TelephonyManager.getSubscriberId()`, which reads the device’s IMSI), is labeled with a taint tag denoting its source. The tag is then propagated along data-flow paths within the code. Any such path that ends up in a release point, or privacy *sink* (e.g. `WebView.loadUrl(...)`, which sends out an HTTP request), triggers a leakage alarm.

The tainting approach effectively reduces leakage judgments to boolean reachability queries. This can potentially lead to false reports, as the real-world example shown in Figure 1 illustrates. This code fragment, extracted from a core library in the Android platform, reads the device’s IMSI number, and then either (i) persists the full number to an error log if the number is invalid (the `loge(...)` call), or (ii) writes a prefix of the IMSI (of length 6) to the standard log while carefully masking away the suffix (of length 9) as ‘x’ characters. Importantly, data flow into the `log(...)` sink is not a privacy problem, because the first 6 digits merely carry model and origin information. Distinctions of this sort are beyond the discriminative power of taint analysis [27].

Quantitative extensions of the core tainting approach have been proposed to address this limitation. A notable example is McCamant and Ernst’s [15] information-flow tracking system, which quantifies flow of secret information by dynamically tracking taint labels at the bit level. Other approaches — based e.g. on distinguishability between secrets [1], the rate of data transmission [14] or the influence inputs have on output values [16] — have also

```
1 String mlmsi = ...; // source
2 // 6 digits <= IMSI (MCC+MNC+MSIN) <= 15 (usually 15)
3 if (mlmsi != null &&
4     (mlmsi.length() < 6 || mlmsi.length() > 15)) {
5     loge(" invalid IMSI:" + mlmsi); // sink
6     mlmsi = null; }
7 log(" IMSI:" + mlmsi.substring(0, 6) + "xxxxxxxxx"); // sink
```

Figure 1: Fragment from an internal Android library, `com.android.internal.telephony.cdma.RuimRecords`, where a prefix of the mobile device’s IMSI number flows into the standard log file

been proposed. While these systems are useful as offline analyses, it is highly unlikely that any of them can be engineered to meet the performance requirements of a real-time monitoring solution due to the high complexity of their underlying algorithms. As an example, McCamant and Ernst report on a workload on which their analysis spent over an hour.

Our Approach We formulate data leakage as a classification problem, which generalizes the source/sink reachability judgment enforced by standard information-flow analysis, permitting richer and more relaxed judgments in the form of statistical classification. The motivating observation is that reasoning about information release is fuzzy in nature. While there are clear examples of legitimate versus illegitimate information release, there are also less obvious cases (e.g., a variant of the example in Figure 1 with a 10- rather than 6-character prefix). A statistical approach, accounting for multiple factors and based on rich data sets, is better able to address these subtleties.

Concretely, we propose Bayesian classification. To label a release point as either legitimate or illegitimate, the Bayesian classifier refers to the “evidence” at that point, and computes the likelihood of each label given the evidence. The evidence consists of feature/value pairs. There are many ways of defining the evidence. In this study, we concentrate on the data arguments flowing into release operations, though we intend to consider other classes of features in the future. (See Section 7.)

Specifically, we induce features over the private values stored on the device, and evaluate these features according to the level of similarity between the private values and those arising at release points. This distinguishes instances where data that is dependent on private values flows into a release point, but its structural and/or quantitative characteristics make it eligible for release, from illegitimate behaviors. Failure to make such distinctions is a common source of false alarms by the tainting approach [4].

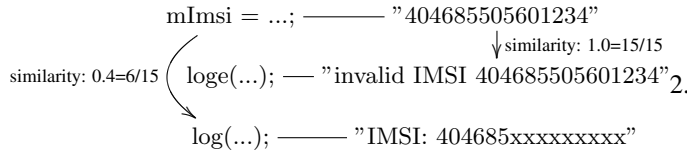


Figure 2: Similarity analysis applied to the code in Figure 1

To illustrate this notion of features, we return to the example in Figure 1. Because the IMSI number is considered private, we define a respective feature *IMSI*. Assume that the concrete IMSI value is “404685505601234”. Then the value arising at the `log(...)` release point is “IMSI: 404685xxxxxxxx”. The quantitative similarity between these two values serves as evidence for the decision whether or not `log(...)` is behaving legitimately. This style of reasoning is depicted in Figure 2.

Evaluation To evaluate our approach, we have implemented the BAYESDROID system for privacy enforcement. [[BAYESDROID is to be featured in the coming version of a commercial product that performs security assessment of web and mobile applications.]] We report on two sets of experiments over BAYESDROID.

First, to measure the accuracy gain thanks to Bayesian analysis, we compared BAYESDROID with the TaintDroid system [4], a highly popular and mature implementation of the tainting approach that is considered both efficient (with average overhead of approximately 10%) and accurate. We applied both BAYESDROID and TaintDroid to the DroidBench suite,¹ which comprises the most mature and comprehensive set of privacy benchmarks currently available. The results suggest dramatic improvement in accuracy thanks to Bayesian elimination of false reports, yielding accuracy scores of 0.96 for BAYESDROID versus 0.66[[0.69]] for TaintDroid.

The second experiment examines the practical value of BAYESDROID by applying it to 54 top-popular mobile apps from Google Play. We evaluate two variants of BAYESDROID, one of which is able to detect a total of 27[[64]] distinct instances of illegitimate information release across 15[[26]] of the applications with only 1 false alarm.

Contributions This paper makes the following principal contributions:

1. Novel approach to leakage detection (Section 2): We present a Bayesian classification alternative to the classic tainting approach. Our approach is more

flexible than taint tracking by permitting statistical weighting of different features as the basis for privacy judgments.

2. Similarity-based reasoning (Section 3): We instantiate the Bayesian approach by applying quantitative similarity judgments over private values and values about to be released. This enables consideration of actual data, rather than only data flow, as evidence for privacy judgments.

3. Implementation and evaluation (Sections 4–5): We have instantiated our approach as the BAYESDROID system, which is about to be featured in a leading security product. We report on two sets of experiments, whose results (i) demonstrate substantial accuracy gain thanks to Bayesian reasoning, and (ii) substantiate the overall effectiveness of BAYESDROID when applied to real-world apps. All the leakage reports by BAYESDROID are publicly available for scrutiny.²

2 The Bayesian Setting

Our starting point is to treat privacy enforcement as a classification problem, being the decision whether or not a given release point is legitimate. The events, or instances, to be classified are (runtime) release points. The labels are *legitimate* and *illegitimate*. Misclassification either yields a false alarm (mistaking benign information release as a privacy threat[[mislabeling a legitimate point as illegitimate]]) or a missed data leak (failing to intercept illegitimate information release[[mislabeling an illegitimate point as legitimate]]).

2.1 Bayes and Naive Bayes

Our general approach is to base the classification on the evidence arising at the release point. Items of evidence may refer to qualitative facts, such as source/sink data-flow reachability, as well as quantitative measures, such as the degree of similarity between private values and values about to be released. These latter criteria are essential in going beyond the question of whether private information is released to also reason about the amount and form of private information about to be released.

A popular classification method, representing this mode of reasoning, is based on Bayes’ theorem (or rule). Given events X and Y , Bayes’ theorem states the following equality:

$$\Pr(Y|X) = \frac{\Pr(X|Y) \cdot \Pr(Y)}{\Pr(X)} \quad (1)$$

where $\Pr(Y|X)$ is the conditional probability of Y given X (i.e., the probability for Y to occur given that X has

¹ <http://sseblog.ec-spride.de/tools/droidbench/>

² <https://www.dropbox.com/sh/ggrcvqsbiubmlb/faSUXmr9xK>

occurred). X is referred to as the *evidence*. Given evidence X , Bayesian classifiers compute the conditional likelihood of each label (in our case, *legitimate* and *illegitimate*).

We begin with the formal background by stating Equation 1 more rigorously. Assume that Y is a discrete-valued random variable, and let $X = [X_1, \dots, X_n]$ be a vector of n discrete or real-valued attributes X_i . Then

$$\Pr(Y = y_k | X_1 \dots X_n) = \frac{\Pr(Y = y_k) \cdot \Pr(X_1 \dots X_n | Y = y_k)}{\sum_j \Pr(Y = y_j) \cdot \Pr(X_1 \dots X_n | Y = y_j)} \quad (2)$$

As Equation 2 hints, training a Bayesian classifier is, in general, impractical. Even in the simple case where the evidence X is a vector of n boolean attributes and Y is boolean, we are still required to estimate a set

$$\theta_{ij} = \Pr(X = x_i | Y = y_j)$$

of parameters, where i assumes 2^n values and j assumes 2 values for a total of $2 \cdot (2^n - 1)$ independent parameters.

Naive Bayes deals with the intractable sample complexity by introducing the assumption of conditional independence, as stated in Definition 2.1 below, which reduces the number of independent parameters sharply to $2n$. Intuitively, conditional independence prescribes that events X and Y are independent given knowledge that event Z has occurred.

Definition 2.1 (Conditional Independence). *Given random variables X , Y and Z , we say that X is conditionally independent of Y given Z iff the probability distribution governing X is independent of the value of Y given Z . That is,*

$$\forall i, j, k. \Pr(X = x_i | Y = y_j, Z = z_k) = \Pr(X = x_i | Z = z_k)$$

Under the assumption of conditional independence, we obtain the following equality:

$$\Pr(X_1 \dots X_n | Y) = \prod_{i=1}^n \Pr(X_i | Y) \quad (3)$$

Therefore,

$$\Pr(Y = y_k | X_1 \dots X_n) = \frac{\Pr(Y = y_k) \cdot \prod_i \Pr(X_i | Y = y_k)}{\sum_j \Pr(Y = y_j) \cdot \prod_i \Pr(X_i | Y = y_j)} \quad (4)$$

2.2 Bayesian Reasoning about Leakage

For leakage detection, conditional independence translates into the requirement that at a release point st , the “weight” of evidence e_1 is not affected by the “weight” of evidence e_2 knowing that st is legitimate/illegitimate. As an example, assuming the evidence is computed as the similarity between private and released values, if st

is known to be a statement sending private data to the network, then the similarity between the IMSI number and respective values about to be released is assumed to be independent of the similarity between location coordinates and respective values about to be released.

The assumption of conditional independence induces a “modular” mode of reasoning, whereby the privacy features comprising the evidence are evaluated independently. This simplifies the problem of classifying a release point according to the Bayesian method into two quantities that we need to clarify and estimate: (i) the likelihood of legitimate/illegitimate release ($\Pr(Y = y_k)$) and (ii) the conditional probabilities $\Pr(X_i | Y = y_k)$.

3 Privacy Features

In this section we develop, based on the mathematical background in Section 2, an algorithm to compute the conditional likelihood of legitimate versus illegitimate data release given privacy features F_i . With such an algorithm in place, given values v_i for the features F_i , we obtain

$$v_{leg} = \Pr(\text{legitimate} | [F_1 = v_1, \dots, F_n = v_n])$$

and

$$v_{illeg} = \Pr(\text{illegitimate} | [F_1 = v_1, \dots, F_n = v_n])$$

Bayesian classification then reduces to comparing between v_{leg} and v_{illeg} , where the label corresponding to the greater of these values is the classification result.

3.1 Feature Extraction

The first challenge that arises is how to define the features (denoted with italicized font: F) corresponding to the private values (denoted with regular font: F). This requires simultaneous consideration of both the actual private value and the “relevant” values arising at the sink statement (or release point). We apply the following computation:

1. Reference value: We refer to the actual private value as the *reference value*, denoting the value of private item F as $\llbracket F \rrbracket$. For the example in Figures 1–2, the reference value, $\llbracket \text{IMSI} \rrbracket$, of the *IMSI* feature would be the device’s IMSI number: $\llbracket \text{IMSI} \rrbracket = “404685505601234”$.
2. Relevant value: We refer to value v about to be released by the sink statement as *relevant* with respect to feature F if there is data-flow connectivity between a source statement reading the value $\llbracket F \rrbracket$ of F and v . Relevant values can thus be computed via information-flow tracking by propagating a unique tag (or label) per each private value, as tools like

TaintDroid already do. Note that for a given feature F , multiple different relevant values may arise at a given sink statement (if the private item F flows into more than one sink argument).

3. Feature value: Finally, given the reference value $\llbracket F \rrbracket$ and a set $\{v_1, \dots, v_k\}$ of relevant values for feature F , the value we assign to F (roughly) reflects the highest degree of pairwise similarity (i.e., minimal distance) between $\llbracket F \rrbracket$ and the values v_i . Formally, we assume a distance metric d . Given d , we define:

$$\llbracket F \rrbracket \equiv \min_{1 \leq i \leq k} \{d(\llbracket F \rrbracket, v_i)\}$$

We leave the distance metric $d(\dots)$ unspecified for now, and return to its instantiation in Section 3.2.

According to our description above, feature values are unbounded in principle, as they represent the distance between the reference value and any data-dependent sink values. In practice, however, assuming (i) the distance metric $d(\dots)$ satisfies $d(x, y) \leq \max\{|x|, |y|\}$, (ii) $\exists c \in \mathbb{N}. |\llbracket F \rrbracket| \leq c$ (as with the IMEI, IMSI, location, etc.), and (iii) $\llbracket F \rrbracket$ is not compared with values larger than it, we can bound $\llbracket F \rrbracket$ by c . In general, any feature can be made finite, with (at most) $n + 1$ possible values, by introducing a privileged “ $\geq n$ ” value, which denotes that the distance between the reference and relevant values is at least n .

3.2 Measuring Distance between Values

To compute a quantitative measure of similarity between data values, we exploit the fact that private data often manifests as strings of ASCII characters [4, 11, 28]. These include e.g. device identifiers (like the IMEI and IMSI numbers), GPS coordinates, inter-application communication (IPC) parameters, etc. This lets us quantify distance between values in terms of string metrics.

Many string metrics have been proposed to date [18]. Two simple and popular metrics, which we have experimented with and satisfy the requirement that $d(x, y) \leq \max\{|x|, |y|\}$, are the following:

Hamming Distance This metric assumes that the strings are of equal length. The Hamming distance between two strings is equal to the number of positions at which the corresponding symbols are different (as indicated by the indicator function $\delta_{c_1 \neq c_2}(\dots)$):

$$\text{ham}(a, b) = \sum_{0 \leq i < |a|} \delta_{c_1 \neq c_2}(a(i), b(i))$$

In another view, Hamming distance measures the number of substitutions required to change one string into the other.

Levenshtein Distance The Levenshtein string metric computes the distance between strings a and b as

Data: Strings u and v
Data: Distance metric d

```

begin
  x ← |u| < |v| ? u : v // min
  y ← |u| ≥ |v| ? u : v // max
  r ← y
  for i = 0 to |y| - |x| do
    y' ← y[i, i + |x| - 1]
    if d(x, y') < r then
      r ← d(x, y')
    end
  end
  end
  return r
end

```

Algorithm 1: The BAYESDROID distance measurement algorithm

$\text{lev}_{a,b}(|a|, |b|)$ (abbreviated as $\text{lev}(|a|, |b|)$), where

$$\text{lev}(i, j) = \begin{cases} \max(i, j) & \text{if } \min(i, j) = 0 \\ \min \begin{pmatrix} \text{lev}(i-1, j) + 1 \\ \text{lev}(i, j-1) + 1 \\ \text{lev}(i-1, j-1) + \delta_{a_i \neq b_j} \end{pmatrix} & \text{otherwise} \end{cases}$$

Informally, $\text{lev}(|a|, |b|)$ is the minimum number of single-character edits — either insertion or deletion or substitution — needed to transform one string into the other. An efficient algorithm for computing the Levenshtein distance is bottom-up dynamic programming [25]. The asymptotic complexity is $O(|a| \cdot |b|)$.

Given string metric $d(x, y)$ and pair (u, v) of reference value u and relevant value v , BAYESDROID computes their distance according to the following steps:

1. BAYESDROID ensures that both u and v are String objects by either (i) invoking `toString()` on reference types or (ii) converting primitive types into Strings (via `String.valueOf(...)`), if the argument is not already of type String.
2. To meet the conservative requirement that $|x| = |y|$ (i.e., x and y are of equal length), BAYESDROID applies Algorithm 1. This algorithm induces a sliding window over the longer of the two strings, whose width is equal to the length of the shorter string. The shorter string is then compared to contiguous segments of the longer string that have the same length. The output is the minimum across all comparisons.

To ensure that comparisons are still meaningful under length adjustment, we decompose private values into indivisible *information units*. These are components of the private value that cannot be broken further, and so comparing them with a shorter value mandates that the shorter value be padded. In our specification, the phone, IMEI and IMSI numbers consist of only one unit of information. The Location object is an example of a data structure that consists of several distinct informa-

tion units. These include the integral and fractional parts of the longitude and latitude values, etc. **BAYESDROID** handles objects that decompose into multiple information units by treating each unit as a separate object and applying the steps above to each of the units in turn. The notion of information units guards BAYESDROID against ill-founded judgments, such as treating release of a single digit from the IMEI number as strong evidence for leakage.

3.3 Estimating Probabilities

The remaining challenge, having clarified what the features are and how their values are computed, is to estimate the probabilities appearing in Equation 4:

- We need to estimate the probability of the *legitimate* event, $\Pr(\text{legitimate})$, where *illegitimate* is the complementary event and thus $\Pr(\text{illegitimate}) = 1 - \Pr(\text{legitimate})$.
- We need to estimate the conditional probabilities $\Pr(F = u | \text{legitimate})$ and $\Pr(F = u | \text{illegitimate})$ for all features F and respective values u .

$\Pr(\text{legitimate})$ can be approximated straightforwardly based on available statistics on the frequency of data leaks in the wild. For the conditional probabilities, assuming feature X_i is discrete valued with j distinct values (per the discussion in Section 3.1 above), we would naively compute the estimated conditional probability θ_{ijk} according to the following equation:

$$\theta_{ijk} = \widehat{\Pr}(X_i = x_{ij} | Y = y_k) = \frac{\#D\{X_i = x_{ij} \wedge Y = y_k\}}{\#D\{Y = y_k\}} \quad (5)$$

The danger, however, is that this equation would produce estimates of zero if the data happens not to contain any training examples satisfying the condition in the numerator. To fix this, we modify Equation 5 as follows:

$$\theta_{ijk} = \widehat{\Pr}(X_i = x_{ij} | Y = y_k) = \frac{\#D\{X_i = x_{ij} \wedge Y = y_k\} + l}{\#D\{Y = y_k\} + l \cdot J} \quad (6)$$

where l is a factor that “smoothen” the estimate by adding in a number of “hallucinated” examples that are assumed to be spread evenly across the possible values of X_i . In Section 5.1, we provide concrete detail on the data sets and parameter values we used for our estimates.

4 The BAYESDROID Algorithm

In this section, we describe the complete BAYESDROID algorithm. We then discuss enhancements of the core algorithm.

4.1 Pseudocode Description

Algorithm 2 summarizes the main steps of BAYESDROID. For simplicity, the description in Algorithm 2

Input: S // privacy specification

```

begin
  while true do
    OnSourceStatement  $r := \text{src } \bar{p}$  :
      // map source to feature
       $f \leftarrow \text{GetFeature } \text{src}$ 
      attach tag  $f$  to  $r$ 
    OnNormalStatement  $r := \text{norm } \bar{p}$  :
      propagate feature tags according to data flow
    OnSinkStatement  $r := \text{snk } \bar{p}$  :
      // map feat.s to param.s with resp. tag
       $\{f \mapsto \bar{p}_f\} \leftarrow \text{ExtractTags } \bar{p}$ 
      foreach  $f \mapsto \bar{p}_f \in \{f \mapsto \bar{p}_f\}$  do
         $u \leftarrow \text{ref } f$ 
         $\delta \leftarrow \min\{d(u, \llbracket p \rrbracket)\}_{p \in \bar{p}_f}$ 
         $f \leftarrow \delta \geq c_f ? \text{“} \geq c_f \text{”} : \delta$ 
      end
      if IsLeakageClassification  $\{f\}$  then
        Alarm snk  $\bar{p}$ 
      end
  end
end

```

Algorithm 2: Outline of the core BAYESDROID algorithm

assumes that source statements serve private data as their return value, though the BAYESDROID implementation also supports other sources (e.g., callbacks like `onLocationChanged(...)`, where the private Location object is passed as a parameter). We also assume that each source maps to a unique privacy feature. Hence, when a source is invoked (i.e., the `OnSourceStatement` event fires), we obtain the unique tag corresponding to its respective feature via the `GetFeature(...)` function. We then attach the tag to the return value r . Normal data flow obeys the standard rules of tag propagation, which are provided e.g. by Enck et al. [4]. (See Table 1 there.)

When an `OnSinkStatement` event is triggered, the arguments flowing into the sink snk are searched for privacy tags, and a mapping from features f to parameters p_f carrying the respective tag is built. The value of f is then computed as the minimal pairwise distance between the parameters $p \in p_f$ and $\text{ref } f$. If this value is greater than some constant c_f defined for f , then the privileged value “ $\geq c_f$ ” is assigned to f . (See Section 3.1.) Finally, the judgment `IsLeakageClassification` is applied over the features whose tags have reached the sink snk . This judgment is executed according to Equation 4.

We illustrate the BAYESDROID algorithm with reference to Figure 3, which demonstrates a real leakage instance in `com.g6677.android.princesshs`, a popular gaming application. In this example, two different private items flow into the sink statement: both the IMEI, read via `getDeviceId()`, and the Android ID, read via

```

1 source : private value
2   TelephonyManager.getDeviceId() : 0000000000000000
3   Settings$Secure.getString (...) : cdf15124ea4c7ad5
4
5 sink : arguments
6   URL.openConnection(...) : app_id=2aec0559c930 ... &
7   android_id=cdf15124ea4c7ad5 \& udid= ... &
8   serial_id = ... & ... &
9   publisher_user_id =0000000000000000

```

Figure 3: True leakage detected by BAYESDROID in com.g6677.android.princesshs

getString(...).

At sink statement `URL.openConnection(...)`, the respective tags *IMEI* and *AndroidID* are extracted. Values are assigned to these features according to the description in Section 3, where we utilize training data, as discussed later in Section 5.1, for Equation 6 (for space, *legitimate/illegitimate* abbreviated as *leg/ilg* and *AndroidID* as *AndID*):

$$\begin{aligned}
\Pr(IMEI \geq 5|leg) &= 0.071 \\
\Pr(IMEI \geq 5|ilg) &= 0.809 \\
\Pr(AndID \geq 5|leg) &= 0.047 \\
\Pr(AndID \geq 5|ilg) &= 0.833
\end{aligned}$$

We then compute Equation 4, where the denominator is the same for both *leg* and *ilg*, and so it suffices to evaluate the nominator (denoted with $\tilde{\Pr}(\dots)$ rather than $\Pr(\dots)$):

$$\begin{aligned}
&\tilde{\Pr}(leg|IMEI \geq 5, AndID \geq 5) = \\
&\Pr(leg) \times \Pr(IMEI \geq 5|leg) \times \Pr(AndID \geq 5|leg) = \\
&\quad 0.66 \times 0.071 \times 0.047 = 0.002 \\
&\tilde{\Pr}(ilg|IMEI \geq 5, AndID \geq 5) = \\
&\Pr(ilg) \times \Pr(IMEI \geq 5|ilg) \times \Pr(AndID \geq 5|ilg) = \\
&\quad 0.33 \times 0.809 \times 0.833 = 0.222
\end{aligned}$$

Our estimates of 0.66 for $\Pr(leg)$ and 0.33 for $\Pr(ilg)$ are again based on training data as explained in Section 5.1. The obtained conditional measure of 0.222 for *ilg* is (far) greater than 0.002 for *leg*, and so BAYESDROID resolves the release instance in Figure 3 as a privacy threat, which is indeed the correct judgment.

4.2 Enhancements

We conclude our description of BAYESDROID by highlighting two extensions of the core algorithm.

```

1 TelephonyManager tm =
2   getSystemService(TELEPHONY_SERVICE);
3 String imei = tm.getDeviceId(); // source
4 String encodedIMEI = Base64Encoder.encode(imei);
5 Log.i(encodedIMEI); // sink

```

Figure 4: Adaptation of the DroidBench Loop1 benchmark, which releases the device ID following Base64 encoding

Beyond Plain Text While many instances of illegitimate information release involve plain text, and can be handled by the machinery in Section 3.1, there are also more challenging scenarios. Two notable challenges are (i) data transformations, whereby data is released following an encoding, encryption or hashing transformation; and (ii) high-volume binary data, such as camera or microphone output. We have extended BAYESDROID to address both of these cases.

We begin with data transformations. As noted earlier, in Section 1, private information is sometimes released following standard hashing/encoding transformations, such as the Base64 scheme. This situation, illustrated in Figure 4, can distort feature values, thereby leading BAYESDROID to erroneous judgments. Fortunately, the transformations that commonly manifest in leakage scenarios are all standard, and there is a small number of such transformations [11].

To account for these transformations, BAYESDROID applies each of them to the value obtained at a source statement, thereby exploding the private value into multiple representations. This is done lazily, once a sink is reached, for performance. This enhancement is specified in pseudocode form in Algorithm 3. The main change is the introduction of a loop that traverses the transformations $\tau \in T$, where the identity transformation, $\lambda x. x$, is included to preserve the (non-transformed) value read at the source. The value assigned to feature f is then the minimum with respect to all transformed values.

Binary data — originating from the microphone, camera or bluetooth adapter — also requires special handling because of the binary versus ASCII representation and, more significantly, its high volume. Our solution is guided by the assumption that such data is largely treated as “uninterpreted” and immutable by application code due to its form and format. This leads to a simple yet effective strategy for similarity measurement, whereby a fixed-length prefix is sampled out of the binary content. Sampling is also applied to sink arguments consisting of binary data.

Heuristic Detection of Relevant Values So far, our description of the BAYESDROID algorithm has relied on

Input: $T \equiv \{\lambda x. x, \tau_1, \dots, \tau_n\} // \text{std. transformations}$

```

begin
  ...
  OnSinkStatement r := snk  $\bar{p}$  :
     $\{f \mapsto \bar{p}_f\} \leftarrow \text{ExtractTags } \bar{p}$ 
    foreach  $f \mapsto \bar{p}_f \in \{f \mapsto \bar{p}_f\}$  do
      foreach  $\tau \in T$  do
         $u \leftarrow \tau(\text{ref } f)$ 
         $\delta \leftarrow \min\{d(u, \llbracket p \rrbracket)\}_{p \in \bar{p}_f}$ 
         $f \leftarrow \min\{\llbracket f \rrbracket, \delta \geq c_f ? \text{"} \geq c_f \text{"} : \delta\}$ 
      end
    end
  ...
end

```

Algorithm 3: BAYESDROID support for standard data transformations

tag propagation to identify relevant values at the sink statement. While this is a robust mechanism to drive feature computation, flowing tags throughout the code also has its costs, incurring runtime overheads of $\geq 10\%$ and affecting the stability of the application due to intrusive instrumentation [4].

These weaknesses of the tainting approach have led us to investigate an alternative method of detecting relevant values. A straightforward relaxation of data-flow tracking is **bounded** (“brute-force”) traversal of the reachable values from the arguments to a sink statement up to some depth bound k : **All values pointed-to by a sink argument or reachable from a sink argument via a sequence of $\leq k$ field dereferences are deemed relevant.** Though in theory this might introduce both false positives (due to irrelevant values that are incidentally similar to the reference value) and false negatives (if k is too small, blocking relevant values from view), in practice both are unlikely, as we confirmed experimentally. (See Section 5.)

For false positives, private values are often unique, and so incidental similarity to irrelevant values is improbable. For false negatives, the arguments flowing into privacy sinks are typically either String objects or simple data structures. Also, because the number of privacy sinks is relatively small, and the number of complex data structures accepted by such sinks is even smaller, it is possible to specify relevant values manually for such data structures. We have encountered only a handful of data structures (e.g., the android.content.Intent class) that motivate a specification.

5 Experimental Evaluation

In this section, we describe the BAYESDROID implementation, and present two sets of experiments that we have conducted to evaluate our approach.

5.1 The BAYESDROID System

Implementation Similarly to existing tools like TaintDroid, BAYESDROID is implemented as an instrumented version of the Android SDK. Specifically, we have instrumented version 4.1.1_r6 of the SDK, which was chosen intentionally to match the most recent version of TaintDroid.³ The experimental data we present indeed utilizes TaintDroid for tag propagation (as required for accurate resolution for relevant values).

Beyond the TaintDroid instrumentation scheme, the BAYESDROID scheme specifies additional behaviors for sources and sinks within the SDK. At source points, a hook is added to record the private value read by the source statement (which acts as a reference value). At sink points, a hook is installed to apply Bayesian reasoning regarding the legitimacy of the sink.

Analogously to TaintDroid, BAYESDROID performs privacy monitoring over APIs for file-system access and manipulation, inter-application and socket communication, reading the phone’s state and location, and sending of text messages. BAYESDROID also monitors the HTTP interface, camera, microphone, bluetooth and contacts. **As explained in Section 4.1, each of the privacy sources monitored by BAYESDROID is mirrored by a tag/feature. The full list of features is as follows: IMEI, IMSI, AndroidID, Location, Microphone, Bluetooth, Camera, Contacts, FileSystem.**

The BAYESDROID implementation is configurable, enabling the user to switch between distance metrics as well as enable/disable information-flow tracking for precise/heuristic determination of relevant values. (See Section 4.2.) In our experiments, we tried both the Levenshtein and the Hamming metrics, but found no observable differences, and so we report the results only once. Our reasoning for why the metrics are indistinguishable is because we apply both to equal-length strings (see Section 3.2), and have made sure to apply the same metric both offline and online, and so both metrics achieve a very similar effect in the Bayesian setting.

Training To instantiate BAYESDROID with the required estimates, as explained in Section 3.3, we applied the following methodology: First, to estimate $\Pr(\text{legitimate})$, we relied on (i) an extensive study by Hornyack et al. spanning 1,100 top-popular free Android apps [11], as well as (ii) a similarly comprehensive study by Enck et al. [5], which also spans a set of 1,100 free apps. According to the data presented in these studies, approximately one out of three release points is illegitimate, and thus $\Pr(\text{legitimate}) = 0.66$ and complementarily $\Pr(\text{illegitimate}) = 1 - 0.66 \approx 0.33$.

³ <http://applanalysis.org/download.html>

For the conditional probabilities $\widehat{\Pr}(X_i = x_{ij} | Y = y_k)$, we queried Google Play for the 100 most popular apps (across all domains) in the geography of one of the authors. We then selected at random 35 of these apps, and analyzed their information-release behavior using debug breakpoints (which we inserted via the adb tool that is distributed as part of the Android SDK).

Illegitimate leaks that we detected offline mainly involved (i) location information and (ii) device and user identifiers, which is consistent with the findings reported by past studies [11, 5]. We confirmed that illegitimate leaks are largely correlated with high similarity between private data and sink arguments, and so we fixed six distance levels for each private item: $[0, 4]$ and ≥ 5 . (See Section 3.1.) Finally, to avoid zero estimates for conditional probabilities while also minimizing data perturbation, we set the “smoothing” factor l in Equation 6 at 1, where the illegitimate flows we detected were in the order of several dozens per private item.

5.2 Experimental Hypotheses

In our experimental evaluation of BAYESDROID, we tested two hypotheses:

1. **H1: Accuracy.** Bayesian reasoning [\[\[about the legitimacy of information release \]\]](#), as implemented in BAYESDROID, yields a significant improvement in [leakage detection](#) accuracy compared to the baseline of information-flow tracking.
2. **H2: Applicability.** For real-life applications, BAYESDROID remains effective under relaxation of the tag-based method for detection of relevant values and its stability improves. [\[\[H2: Power of Bayesian Analysis. BAYESDROID remains effective under relaxation of the tag-based method for detection of relevant values, and its stability and applicability to real-world applications improve. \]\]](#)

5.3 H1: Accuracy [\[\[of BAYESDROID \]\]](#)

To assess the accuracy of BAYESDROID, we compared it to that of TaintDroid — a state-of-the-art information-flow tracking tool for privacy enforcement. Our experimental settings and results are described below.

Subjects We applied both TaintDroid and BAYESDROID to DroidBench — an independent and publicly available collection of benchmarks serving as testing ground for both static and dynamic privacy enforcement algorithms. DroidBench models a large set of realistic challenges in leakage detection, including precise tracking of sensitive data through containers, handling of call-

backs, field and object sensitivity, lifecycle modeling, inter-app communication, reflection and implicit flows.

The DroidBench suite consists of 50 cases. We excluded from our experiment (i) 8 benchmarks that crash at startup, as well as (ii) 5 benchmarks that leak data via callbacks that we did not manage to trigger (e.g., `onLowMemory()`), as both TaintDroid and BAYESDROID were naturally unable to detect leakages in these two cases. The complete list of benchmarks that we used can be found in Table 4 of Appendix B.

Methodology For each benchmark, we measured the number of true-positive (TP), false-positive (FP) and false-negative (FN) results, and then calculated the *precision* and *recall* of each tool using the formulas below:

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

High precision implies that a technique returns few irrelevant results, whereas high recall implies that it misses only few relevant ones. Since ideal techniques have both high recall and high precision, the F-measure is commonly used to combine both precision and recall into a single measure. The F-measure is defined as the harmonic mean of precision and recall, and is calculated as follows:

$$F\text{-Measure} = 2 * \frac{Precision * Recall}{Precision + Recall}$$

The value of F-measure is high only when both precision and recall are high. We thus use the F-measure to evaluate the accuracy of each technique.

Results The results obtained for both TaintDroid and BAYESDROID on version 1.1 of DroidBench are summarized in Table 1 and presented in detail in Table 4. The findings reported by BAYESDROID are also publicly available.⁴

Overall, TaintDroid detects 31 true leakages while also reporting 17 false positives, whereas our approach suffers from 2 false negatives, discovering 29 of the true leakages while flagging only 1 false positive. We averaged the F-measure calculation across all benchmarks (0.66 for TaintDroid vs. 0.96 for BAYESDROID). We also provide the “global” score computed by considering all TP, FP and FN results together in parentheses (0.78 vs 0.95).

The results mark BAYESDROID as visibly more accurate than TaintDroid. To further confirm this observation,

⁴ See archive file `droidbench.zip` available online at <https://www.dropbox.com/sh/ggrcvqsbiubmlb/faSUXmr9xK>.

	TPs	FPs	FNs	Precision	Recall	F-measure
TaintDroid	31	17	0	0.66	1.00	0.66 (0.78)
BAYESDROID	29	1	2	0.98	0.97	0.96 (0.95)

Table 1: Accuracy of BAYESDROID and TaintDroid on DroidBench

we performed a Welch’s T-Test — an adaptation of Student’s T-Test intended to test the null hypothesis that the two population means (in our case, average F-measures) are equal. Welch’s variant of the T-Test is used when two samples are having possibly unequal variances [26]. The result of applying the test rejects the null hypothesis and demonstrates that BAYESDROID performs significantly better than TaintDroid ($p\text{-value} \ll 0.05$), which *confirms H1*.

Discussion Analysis of the per-benchmark findings reveals the following: First, the 2 false negatives of BAYESDROID on ImplicitFlow1 are both due to proprietary (i.e., non-standard) data transformations, which are outside the current scope of BAYESDROID. An illustrative fragment from the ImplicitFlow1 code is shown in Figure 5. The obfuscateIMEI(...) transformation maps IMEI digits to English letters, which is a non-standard behavior that is unlikely to arise in an authentic app.

The true positive reported by BAYESDROID, in common with TaintDroid, is on release of sensitive data to the file system, albeit using the MODE_PRIVATE flag, which does not constitute a leakage problem in itself. This can be resolved by performing Bayesian reasoning not only over argument values, but also over properties of the sink API (in this case, the storage location mapped to a file handle). We intend to implement this enhancement.

Beyond the false alarm in common with BAYESDROID, TaintDroid has multiple other sources of imprecision. The main reasons for its false positives are

- coarse modeling of containers, mapping their entire contents to a single taint bit, which accounts e.g. for the false alarms on `ArrayAccess{1,2}` and `HashMapAccess1`;
- field and object insensitivity, resulting in false alarms on `FieldSensitivity{2,4}` and `ObjectSensitivity{1,2}`; and more fundamentally,
- ignoring of data values, which causes TaintDroid to issue false warnings on `LocationLeak{1,2}` even when location reads fail, yielding a `Location` object without any meaningful information.

The fundamental reason for these imprecisions is to constrain the overhead of TaintDroid, such that it can meet the performance demands of online privacy enforcement.

```

1 TelephonyManager tm =
2   getSystemService(TELEPHONY_SERVICE);
3 String imei = tm.getDeviceId(); //source
4 String obfuscatedIMEI = obfuscateIMEI(imei); ...;
5 Log.i(imei); // sink
6
7 private String obfuscateIMEI(String imei) {
8   String result = "";
9   for (char c : imei.toCharArray()) {
10    switch(c) {
11      case '0': result += 'a'; break;
12      case '1': result += 'b'; break;
13      case '2': result += 'c'; break;
14      ...; } }

```

Figure 5: Fragment from the DroidBench ImplicitFlow1 benchmark, which applies a proprietary transformation to private data

BAYESDROID is able to accommodate such optimizations while still ensuring high accuracy.

5.4 H2: Applicability

The second aspect of the evaluation compared between two versions of BAYESDROID, whose sole difference lies in the method used for detecting relevant values: In one configuration (T-BD), relevant values are detected via tag propagation. The other configuration (H-BD) uses the heuristic detailed in Section 4.2 of treating all values reachable from sink arguments (either directly or via the heap graph) up to a depth bound of k as relevant, which places more responsibility on Bayesian reasoning. We set k at 3 based on manual review of the data structures flowing into privacy sinks.

We designed a parametric benchmark application to quantify the overhead reduction imposed by the H-BD variant of BAYESDROID. The application[[First, to quantify overhead reduction, we designed a parametric benchmark application, which]] consists of a simple loop that flows the device IMEI into a log file. Loop iterations perform intermediate data propagation steps. We then performed a series of experiments — over the range of 1 to 19 propagation steps — to quantify the relative

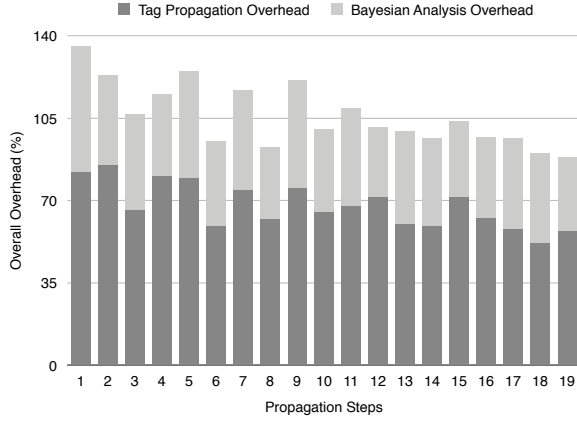


Figure 6: Overhead breakdown into tag propagation (dark grey) and Bayesian analysis at sink (pale grey), where horizontal (X) axis denotes the number of propagation steps and vertical (Y) axis denotes overhead

overhead of tag propagation versus Bayesian analysis.

The results, presented in Figure 6, suggest that the overhead of tag propagation is more dominant than that of Bayesian analysis (with a ratio of roughly 2:1), even when the set of relevant values is naively over approximated. Discussion of the methodology underlying this experiment is provided in Appendix A.

Yet, in the general case,[[In summary]] H-BD trades overhead reduction for accuracy. H2 then asserts that, *in practice*, the tradeoff posed by H-BD is effective. Below, we discuss our empirical evaluation of this hypothesis.

Subjects For accuracy comparison, we selected a set of real-world applications using the following methodology to avoid evaluators’ bias: We started from the 65 Google Play apps not chosen for the training phase. We then excluded 8 apps that do not have permission to access sensitive data and/or perform release operations (i.e., their manifest does not declare sufficient permissions out of INTERNET, READ_PHONE_STATE, SEND_SMS, etc), as well as 3 apps that we did not manage to install properly, resulting in 54 apps that installed successfully and exercise privacy sources and sinks.

The complete list of the application we used is given in Table 5 of Appendix B. A subset of the applications, for which at least one leakage was detected, is also listed in Table 3.

Methodology We deployed the apps under the two BAYESDROID configurations. Each execution was done from a clean starting state. The third column of both Tables 3 and 5 denotes whether our exploration of the app was exhaustive. By that we mean exercising all the

UI points exposed by the app in a sensible order. Ideally we would do so for all apps. However, (i) some of the apps, and in particular gaming apps, had stability issues, and (ii) certain apps require SMS-validated sign in, which we did not perform. We did, however, create Facebook, Gmail and Dropbox accounts to log into apps that demand such information yet do not ask for SMS validation. We were also careful to execute the exact same crawling scenario under both the T-BD and H-BD configurations. We comment, from our experience, that most data leaks happen when an app launches, and initializes advertising/analytics functionality, and so for apps for which deep crawling was not possible the results are still largely meaningful.

For comparability between the H-BD and T-BD configurations, we counted different dynamic reports involving the same pair of source/sink APIs as a single leakage instance. We manually classified the findings into true positives and false positives. For this classification, we scrutinized the reports by the two configurations, and also — in cases of uncertainty — decompiled and/or reran the app to examine its behavior more closely. As in the experiment described in Section 5.3, we then calculated the precision, recall and F-measure for each of the tools.

Results The results obtained for H-BD and T-BD are summarized in Table 2. Table 3 summarizes the findings reported by both H-BD and T-BD at the granularity of privacy items: the device number, identifier and location, while Table 5 provides a detailed description of the results across all benchmarks including those on which no leakages were detected. The warnings reported by the H-BD configuration are also publicly available for review.⁵

As Table 2 indicates the H-BD variant is more accurate than the T-BD variant overall (F-measure of 0.99 vs. 0.84). As in the experiment described in Section 5.3, we performed a Welch’s T-Test on the F-Measure results to further confirm this observation. The test establishes that the difference between the accuracy of the tools is significant ($p\text{-value}=0.003 < 0.05$). Moreover, H-BD has a lower number of crashes and lower runtime overhead, which confirms H2.

Discussion To give the reader a taste of the findings, we present in Figures 7–8 two examples of potential leakages that BAYESDROID (both the H-BD and the T-BD configurations) deemed legitimate. The instance in Figure 7 reflects the common scenario of obtaining the current (or last known) location, converting it

⁵ See archive file [realworldapps.zip](https://www.dropbox.com/sh/ggrcvqsbkiubmlb/faSUXmr9xK) available online at <https://www.dropbox.com/sh/ggrcvqsbkiubmlb/faSUXmr9xK>.

	TPs	FPs	FNs	Precision	Recall	F-measure	Crashes
H-BD	27	1	0	0.99	1.00	0.99 (0.98)	12
T-BD	13	0	9	1.00	0.84	0.84 (0.74)	22

Table 2: Accuracy of H-BD and T-BD BAYESDROID configurations

App	Domain	Deep crawl?	H-BD			T-BD		
			no.	dev. ID	loc.	no.	dev. ID	loc.
atsoft.games.smgame	games/arcade	✓		✓	✓		✓	✓
com.antivirus	communication	✓		✓			✓	
com.appershopper.ios7lockscreen	personalization		✓	✓	✓			✓
com.bestcoolfungames.antsmasher	games/arcade	✓			✓			✓
com.bitfitlabs.fingerprint.lockscreen	games/casual			✓				
com.cleanmaster.mguard	tools	✓		✓			✓	
com.coolfish.cathairsalon	games/casual	✓		✓				
com.coolfish.snipershooting	games/action	✓		✓				
com.digisoft.TransparentScreen	entertainment	✓			✓			✓
com.g6677.android.cbaby	games/casual			✓				
com.g6677.android.chospital	games/casual			✓				
com.g6677.android.design	games/casual			✓				
com.g6677.android.pnailspa	games/casual			✓				
com.g6677.android.princesshs	games/casual			✓				
com.goldtouch.mako	news	✓		✓			✓	
15		8	1	13	4	0	4	4

Table 3: Findings due to the H-BD and T-BD BAYESDROID configurations on 15/54 top-popular mobile apps (phone number abbreviated as no.; device ID as dev. ID; and location as loc.)

```

1 source : private value
2   GeoCoder.getFromLocation(...) : [ Lat: ..., Long: ...,
3     Alt: ..., Bearing: ..., ..., IL ]
4
5 sink : arguments
6   WebView.loadUrl(...) : http://linux.appwiz.com/
7     profile /72/72_exitad.html?
8     p1=RnVsbCtBbmRyb2lkK29uK0VtdWxhdG9y&
9     p2=Y2RmMTUxMjRIYTRjN2FkNQ%3d%3d&
10    ...
11    LOCATION=IL&
12    ...
13    MOBILE_COUNTRY_CODE=&
14    NETWORK=WIFI

```

Figure 7: Suppressed spurious leakage on ios7lockscreen

```

1 source : private value
2   Settings$Secure.getString (...) : cdf15124ea4c7ad5
3
4 sink : arguments
5   FileOutputStream.write (...) :
6     <?xml version='1.0' encoding='utf-8'
7     standalone='yes'
8     ?><map><string
9     name="openudid">cdf15124ea4c

```

Figure 8: Suppressed spurious leakage on fruitnin-jafree

into one or more addresses, and then releasing only the country or zip code. In the second instance, in Figure 8, the 64-bit Android ID — generated when the user first sets up the device — is read via a call to `Settings$Secure.getString(ANDROID_ID)`. At the release point, into the file system, only a prefix of the Android ID consisting of the first 12 digits is published.

As indicated in Table 3, the findings by H-BD are more complete: it detects 27 leakages (vs. 13 for T-BD), with no false negative results and only one false positive. We attribute that to (i) the intrusive instrumentation required for tag propagation, which can cause instabilities, and (ii) inability to track tags through native code, as discussed below.

The T-BD variant introduces significantly more instability than the H-BD variant, causing illegal application behaviors in 21 cases compared to only 12 under H-BD. We have investigated this large gap between the H-BD and T-BD configurations, including by decompiling the subject apps. Our analysis links the vast majority of illegal behaviors to limitations that TaintDroid casts on loading of third-party libraries. For this reason, certain functionality is not executed, also leading to exceptional app states, which both inhibit certain data leaks.⁶

The secondary reason why H-BD is able to detect more leakages, e.g., in `com.bitfitlabs.fingerprint.lockscreen`, is that this benchmark makes use of the `mobileCore` module,⁷ which is a highly optimized and obfuscated library. We suspect that data flow breaks within this library, though we could not fully confirm this.

At the same time, the loss in accuracy due to heuristic identification of relevant values is negligible, as suggested by the discussion in Section 4.2. H-BD triggers only one false alarm, on `iso7lockscreen`, which is due to overlap between irrelevant values: extra information on the `Location` object returned by a call to `LocationManager.getLastKnownLocation(...)` and unrelated meta-data passed into a `ContextWrapper.startService(...)` request. Finally, as expected, H-BD does not incur false negatives.

6 Related Work

As most of the research on privacy monitoring builds on the tainting approach, we survey related research mainly in this space. We also mention several specific studies in other areas.

⁶ For a technical explanation, see forum comment by William Enck, the TaintDroid moderator, at <https://groups.google.com/forum/#!topic/android-security-discuss/U1fteEX26bk>.

⁷ <https://www.mobilecore.com/sdk/>

Realtime Techniques The state-of-the-art system for realtime privacy monitoring is TaintDroid [4]. TaintDroid features tolerable runtime overhead of about 10%, and can track taint flow not only through variables and methods but also through files and messages passed between apps. TaintDroid has been used, extended and customized by several follow-up research projects. Jung et al. [12] enhance TaintDroid to track additional sources (including contacts, camera, microphone, etc). They used the enhanced version in a field study, which revealed 129 of the 223 apps they studied as vulnerable. 30 out of 257 alarms were judged as false positives. The Kynoid system [21] extends TaintDroid with user-defined security policies, which include e.g. temporal constraints on data processing as well as restrictions on destinations to which data is released.

The main difference between BAYESDROID and the approaches above, which all apply information-flow tracking, is that BAYESDROID exercises “fuzzy” reasoning, in the form of statistical classification, rather than enforcing a clear-cut criterion. As part of this, BAYESDROID factors into the privacy judgment the data values flowing into the sink statement, which provides additional evidence beyond data flow.

Quantitative Approaches Different approaches have been proposed for quantitative information-flow analysis, all unified by the observation that data leakage is a quantitative rather than boolean judgment. McCamant and Ernst [15] present an offline dynamic analysis that measures the amount of secret information that can be inferred from a program’s outputs, where the text of the program is considered public. Their approach relies on taint analysis at the bit level. It provides a sound upper bound on the actual information flow, but only with respect to the observed runs. Newsome et al. [16] develop complementary techniques to bound a program’s *channel capacity* using decision procedures (SAT and #SAT solvers). They apply these techniques to the problem of false positives in dynamic taint analysis. Backes et al. [1] measure leakage in terms of indistinguishability, or equivalence, between outputs due to different secret artifacts. Their characterization of equivalence relations builds on the information-theoretic notion of entropy. Budi et al. [2] propose *kb*-anonymity, a model inspired by *k*-anonymity for ensuring safe release of private data for the purposes of testing and debugging. Like *k*-anonymity, *kb*-anonymity replaces certain information in the original data for privacy preservation, but beyond that, it also ensures that the replaced data does not lead to divergent program behaviors, and thus testing and debugging are still effective.

While these proposals have all been shown useful, none of these approaches has been shown to be efficient

enough to meet realtime constraints. The algorithmic complexity of computing the information-theoretic measures introduced by these works seriously limits their applicability in a realtime setting. Our approach, instead, enables a quantitative/probabilistic mode of reasoning that is simultaneously lightweight, and therefore acceptable for online monitoring, by focusing on relevant features that are efficiently computable.

Static and Hybrid Analysis The PiOS tool [3] performs static leakage analysis of iOS applications. PiOS combines backward slicing with forward constant propagation for accurate resolution of method calls. It then performs standard taint analysis atop the resulting call graph, where path length is (unsoundly) limited to a maximum of 100 basic blocks. The authors report that over half of the 1,400 apps they analyzed leak the unique ID of their host device. Graa et al. [7] propose a hybrid tainting technique that augments dynamic taint analysis with certain classes of implicit flows (namely, those due to conditional statements). To do so, they apply intraprocedural static analysis to the control flow graph of a method at load time to detect variable assignments that occur only under one conditional branch. The FlowDroid tool [6] performs static leakage analysis of Android applications using tainting. It accounts for the Android lifecycle model, handles callbacks and features flow, field and object sensitivity. FlowDroid builds on, and extends, previous taint analyses designed for the server side [24, 23] as well as client side [9] of web applications.

In our experience, even precise static analysis techniques often suffer from a large number of false findings. These are due both to the qualitative (and coarse) nature of taint analysis and to other approximations (e.g., call-graph imprecisions). As for hybrid techniques, current techniques address challenges like implicit flows, where our experience suggests that for effective detection of data leakage by authentic apps, a more important dimension is the form and quantity of released information. In Section 7, we outline a hybrid variant of BAYESDROID that we plan to develop.

Offline Algorithms The AppsPlayground framework [19], built atop TaintDroid, performs dynamic analysis of mobile apps. It features automatic triggering of system events, including usage of contextual information to provide meaningful textual input. The AdRisk system [8] identifies potential privacy and security risks due to embedded in-app advertisement libraries. It first decouples the embedded ad libraries from the host app, and then applies state analysis to each of the libraries to approximate its behavior. AdSplit [22] handles third-party advertising services differently by recompiling the

app to extract such services, which consequently run as separate processes under separate user IDs.

Our approach is extensible to the offline setting, e.g. by piggybacking on unit/integration tests and detecting leakage bugs during the development and testing phases of a mobile app. Our technique can also be integrated with frameworks like AppsPlayground, which provide automated crawling capabilities, to perform offline testing for leakage vulnerabilities.

7 Conclusion and Future Work

In this paper, we articulated the problem of privacy enforcement in mobile systems as a classification problem. The traditional approach of information-flow tracking then becomes a “binary” classifier, which equates data leakage with source/sink data flow. We explored an alternative approach based on statistical classification, which addresses more effectively the inherent fuzziness in leakage judgements, accounting beyond data flow also for the amount and form of private information about to be released. We have instantiated our approach as the BAYESDROID system, which is about to be featured in a commercial security product. Our experimental data establishes the high accuracy of BAYESDROID as well as its applicability to real-world mobile apps.

Moving forward, we have two main objectives. The first is to extend BAYESDROID with additional feature types. Specifically, we would like to account for (i) sink properties, such as file access modes (private vs public), the target URL of HTTP communication (same domain or third party), etc; as well as (ii) the history of privacy-relevant API invocations up to the release point (checking e.g. if/which declassification operations were invoked). Our second objective is to optimize our flow-based method for detecting relevant values (see Section 3.1) by applying (offline) static taint analysis to the subject program, e.g. using the FlowDroid tool [6].

References

- [1] M. Backes, B. Kopf, and A. Rybalchenko. Automatic discovery and quantification of information leaks. In *Proceedings of the 2009 30th IEEE Symposium on Security and Privacy*, pages 141–153, 2009.
- [2] A. Budi, D. Lo, L. Jiang, and Lucia. kb-anonymity: a model for anonymized behaviour-preserving test and debugging data. In *PLDI*, pages 447–457, 2011.

- [3] M. Egele, C. Kruegel, E. Kirda, and G. Vigna. PiOS: Detecting Privacy Leaks in iOS Applications. In *NDSS*, 2011.
- [4] W. Enck, P. Gilbert, B. Chun, L. P. Cox, J. Jung, P. McDaniel, and A. N. Sheth. Taintdroid: an information-flow tracking system for realtime privacy monitoring on smartphones. In *Proceedings of the 9th USENIX conference on Operating systems design and implementation*, pages 1–6, 2010.
- [5] W. Enck, D. Ocateau, P. McDaniel, and S. Chaudhuri. A study of android application security. In *Proceedings of the 20th USENIX Conference on Security*, pages 21–21, 2011.
- [6] C. Fritz, S. Arzt, S. Rasthofer, E. Bodden, A. Bartel, J. Klein, Y. Traon, D. Ocateau, and P. McDaniel. Flowdroid: Precise context, flow, field, object-sensitive and lifecycle-aware taint analysis for android apps. Technical Report TUD-CS-2013-0113, Secure Software Engineering Group, EC Spride, 2013.
- [7] M. Graa, N. Cuppens-Boulahia, F. Cuppens, and A. Cavalli. Detecting control flow in smartphones: combining static and dynamic analyses. In *Proceedings of the 4th international conference on Cyberspace Safety and Security*, pages 33–47, 2012.
- [8] M. C. Grace, W. Zhou, X. Jiang, and A. Sadeghi. Unsafe exposure analysis of mobile in-app advertisements. In *Proceedings of the fifth ACM conference on Security and Privacy in Wireless and Mobile Networks*, pages 101–112, 2012.
- [9] S. Guarnieri, M. Pistoia, O. Tripp, J. Dolby, S. Teihet, and R. Berg. Saving the world wide web from vulnerable javascript. In *ISSTA*, pages 177–187, 2011.
- [10] S. Holavanalli, D. Manuel, V. Nanjundaswamy, B. Rosenberg, F. Shen, S. Y. Ko, and L. Ziarek. Flow permissions for android. In *ASE*, pages 652–657, 2013.
- [11] P. Hornyack, S. Han, J. Jung, S. E. Schechter, and D. Wetherall. These aren’t the droids you’re looking for: retrofitting android to protect data from imperious applications. In *ACM Conference on Computer and Communications Security*, pages 639–652, 2011.
- [12] J. Jung, S. Han, and D. Wetherall. Short paper: enhancing mobile application permissions with runtime feedback and constraints. In *Proceedings of the second ACM workshop on Security and privacy in smartphones and mobile devices*, SPSM ’12, pages 45–50, 2012.
- [13] B. Livshits and J. Jung. Automatic mediation of privacy-sensitive resource access in smartphone applications. In *Proceedings of the 22Nd USENIX Conference on Security*, pages 113–130, 2013.
- [14] G. Lowe. Quantifying information flow. In *Proceedings of the 15th IEEE workshop on Computer Security Foundations*, CSFW ’02, pages 18–31, 2002.
- [15] S. McCamant and M. D. Ernst. Quantitative information flow as network flow capacity. In *PLDI*, pages 193–205, 2008.
- [16] J. Newsome, S. McCamant, and D. Song. Measuring channel capacity to distinguish undue influence. In *Proceedings of the ACM SIGPLAN Fourth Workshop on Programming Languages and Analysis for Security*, pages 73–85, 2009.
- [17] J. Newsome and D. X. Song. Dynamic taint analysis for automatic detection, analysis, and signature generation of exploits on commodity software. In *NDSS*, 2005.
- [18] J. Piskorski and M. Sydow. String distance metrics for reference matching and search query correction. In *Proceedings of the 10th International Conference on Business Information Systems*, pages 353–365, 2007.
- [19] V. Rastogi, Y. Chen, and W. Enck. Appsplayground: automatic security analysis of smartphone applications. In *Proceedings of the third ACM conference on Data and application security and privacy*, pages 209–220, 2013.
- [20] G. Sarwar, O. Mehani, R. Boreli, and M. A. Kaffar. On the effectiveness of dynamic taint analysis for protecting against private information leaks on android-based devices. In *SECRYPT*, pages 461–468, 2013.
- [21] D. Schreckling, J. Posegga, J. Köstler, and M. Schaff. Kynoid: real-time enforcement of fine-grained, user-defined, and data-centric security policies for android. In *WISTP*, pages 208–223, 2012.
- [22] S. Shekhar, M. Dietz, and D. S. Wallach. Ad-split: Separating smartphone advertising from applications. In *Proceedings of the 21st USENIX Security Symposium*, pages 553–567, 2012.

- [23] O. Tripp, M. Pistoia, P. Cousot, R. Cousot, and S. Guarnieri. Andromeda: Accurate and scalable security analysis of web applications. In *FASE*, pages 210–225, 2013.
- [24] O. Tripp, M. Pistoia, S. J. Fink, M. Sridharan, and O. Weisman. Taj: effective taint analysis of web applications. pages 87–97, 2009.
- [25] R. A. Wagner and M. J. Fischer. The string-to-string correction problem. *J. ACM*, 21(1):168–173, 1974.
- [26] Bernard L Welch. The generalization of student’s problem when several different population variances are involved. *Biometrika*, 34(1–2):28–35, 1947.
- [27] D. Wetherall, D. Choffnes, B. Greenstein, S. Han, P. Hornyack, J. Jung, S. Schechter, and X. Wang. Privacy revelations for web and mobile apps. In *Proceedings of the 13th USENIX conference on Hot topics in operating systems*, pages 21–21, 2011.
- [28] Z. Yang, M. Yang, Y. Zhang, G. Gu, P. Ning, and X. S. Wang. Appintent: analyzing sensitive data transmission in android for privacy leakage detection. In *Proceedings of the 2013 ACM SIGSAC Conference on Computer and Communications Security*, pages 1043–1054, 2013.

A Overhead Measurement: Methodology

To complete the description in Section 5.4, we now detail the methodology governing our overhead measurements. The behavior of the benchmark app is governed by two user-controlled values: (i) the length ℓ of the source/sink data-flow path (which is proportional to the number of loop iterations) and (ii) the number m of values reachable from sink arguments.

Based on our actual benchmarks, as well as data reported in past studies [24], we defined the ranges $1 \leq \ell \leq 19$ and $1 \leq m \leq 13 = \sum_{n=0}^2 3^n$. We then ran the parametric app atop a “hybrid” configuration of BAYESDROID that simultaneously propagates tags and treats all the values flowing into a sink as relevant.

For each value of ℓ , we executed the app 51 times, picking a value from the range $[0, 2]$ for n uniformly at random in each of the 51 runs. We then computed the average overhead over the runs, excluding the first (cold) run to remove unrelated initialization costs. The stacked columns in Figure 6 each correspond to a unique value of ℓ .

B Detailed Results [[on Real-world Apps]]

Table 4 summarizes the results of the H1 experiment described in Section 5.3. For each of the benchmarks, it specifies the number of true positive, false positive and false negative findings, as well as the calculated precision, recall and F-measure for the compared tools – BAYESDROID and TaintDroid. The benchmarks for which the tools differ are highlighted, for convenience. The F-measure is calculated by (a) averaging the corresponding values across all benchmarks and (b) directly by considering the total number of TPs, FPs and FNs results. Both values are given in Table 4, outside and inside of the parentheses, respectively: 0.66 for TaintDroid vs. 0.96 for BAYESDROID in the former case, and 0.78 vs 0.95 in the latter.

Similarly, Table 5 summarizes the results of the H2 experiment described in Section 5.4. The first two columns of Table 5 list the applications and their respective domain, and the third column denotes whether crawling was exhaustive. Then, the number of crashes, true positive, false positive and false negative findings, as well as the calculated precision, recall and F-measure are reported for both the H-BD and T-BD variant of BAYESDROID. As in Table 4, the F-measure is calculated by (a) averaging the corresponding value across all benchmarks (0.99 for H-BD vs. 0.84 for T-BD), and (b) directly by considering the total number of TPs, FPs and FNs (0.98 vs. 0.74, in parentheses).

Benchmark	BAYESDROID						TaintDroid					
	TPs	FPs	FNs	Precision	Recall	F-measure	TPs	FPs	FNs	Precision	Recall	F-measure
ActivityCommunication1	1	0	0	1	1	1	1	0	0	1	1	1
ActivityLifecycle1	1	0	0	1	1	1	1	0	0	1	1	1
ActivityLifecycle2	1	0	0	1	1	1	1	0	0	1	1	1
ActivityLifecycle4	1	0	0	1	1	1	1	0	0	1	1	1
Library2	1	0	0	1	1	1	1	0	0	1	1	1
Obfuscation1	1	0	0	1	1	1	1	0	0	1	1	1
PrivateDataLeak3	1	1	0	0.5	1	0.667	1	1	0	0.5	1	0.667
AnonymousClass1	0	0	0	1	1	1	0	1	0	0	1	0
ArrayAccess1	0	0	0	1	1	1	0	1	0	0	1	0
ArrayAccess2	0	0	0	1	1	1	0	1	0	0	1	0
HashMapAccess1	0	0	0	1	1	1	0	1	0	0	1	0
Button1	1	0	0	1	1	1	1	0	0	1	1	1
Button3	2	0	0	1	1	1	2	0	0	1	1	1
Ordering1	0	0	0	1	1	1	0	2	0	0	1	0
RegisterGlobal1	1	0	0	1	1	1	1	0	0	1	1	1
DirectLeak1	1	0	0	1	1	1	1	0	0	1	1	1
FieldSensitivity2	0	0	0	1	1	1	0	1	0	0	1	0
FieldSensitivity3	1	0	0	1	1	1	1	0	0	1	1	1
FieldSensitivity4	0	0	0	1	1	1	0	1	0	0	1	0
ImplicitFlow1	0	0	2	1	0	0	2	0	0	1	1	1
InheritedObjects1	1	0	0	1	1	1	1	0	0	1	1	1
ListAccess1	0	0	0	1	1	1	0	1	0	0	1	0
LocationLeak1	0	0	0	1	1	1	0	2	0	0	1	0
LocationLeak2	0	0	0	1	1	1	0	2	0	0	1	0
Loop1	1	0	0	1	1	1	1	0	0	1	1	1
Loop2	1	0	0	1	1	1	1	0	0	1	1	1
ApplicationLifecycle1	1	0	0	1	1	1	1	0	0	1	1	1
ApplicationLifecycle3	1	0	0	1	1	1	1	0	0	1	1	1
MethodOverride1	1	0	0	1	1	1	1	0	0	1	1	1
ObjectSensitivity1	0	0	0	1	1	1	0	1	0	0	1	0
ObjectSensitivity2	0	0	0	1	1	1	0	2	0	0	1	0
Reflection1	1	0	0	1	1	1	1	0	0	1	1	1
Reflection2	1	0	0	1	1	1	1	0	0	1	1	1
Reflection3	1	0	0	1	1	1	1	0	0	1	1	1
Reflection4	1	0	0	1	1	1	1	0	0	1	1	1
SourceCodeSpecific1	5	0	0	1	1	1	5	0	0	1	1	1
StaticInitialization1	1	0	0	1	1	1	1	0	0	1	1	1
Total	29	1	2				31	17	0			
Average				0.98	0.97	0.96 (0.95)				0.66	1.00	0.66 (0.78)

Table 4: Detailed summary of the results of the H1 experiment described in Section 5.3

App	Domain	Deep crawl?	H-BD							T-BD						
			Crashes	TPs	FPs	FNs	Precision	Recall	F-measure	Crashes	TPs	FPs	FNs	Precision	Recall	F-measure
air.au.com.metro.DumbWaysToDie	games/casual		✓	0	0	0	1	1	1	✓	0	0	0	1	1	1
at.nerbrothers.SuperJump	games/arcade			0	0	0	1	1	1	✓	0	0	0	1	1	1
atsoft.games.smggame	games/arcade	✓		4	0	0	1	1	1		4	0	0	1	1	1
com.antivirus	communication	✓		1	0	0	1	1	1		1	0	0	1	1	1
com.appersshopper.ios7.lockscreen	personalization			5	1	0	0.833	1	0.909		2	0	1	1	0.667	0.8
com.applicaster.il.hotvod	entertainment	✓		0	0	0	1	1	1		0	0	0	1	1	1
com.appstar.callrecorder	tools			0	0	0	1	1	1		0	0	0	1	1	1
com.awesomecargames.mountainclimbrace.1	games/racing		✓	0	0	0	1	1	1	✓	0	0	0	1	1	1
com.bestcoolfungames.antsmasher	games/arcade	✓		2	0	0	1	1	1		2	0	0	1	1	1
com.bigduckgames.flow	games/puzzles			0	0	0	1	1	1	✓	0	0	0	1	1	1
com.bitfitlabs.fingerprint.lockscreen	games/casual			2	0	0	1	1	1		0	0	1	1	0	0
com.channel2.mobile.ui	news	✓		0	0	0	1	1	1		0	0	0	1	1	1
com.chillingo.parkingmaniafree.android.rowgplay	games/racing		✓	0	0	0	1	1	1	✓	0	0	0	1	1	1
com.cleanmaster.mguard	tools	✓		1	0	0	1	1	1		1	0	0	1	1	1
com.coolfish.cathairsalon	games/casual	✓		2	0	0	1	1	1	✓	0	0	1	1	0	0
com.coolfish.snipershooting	games/action	✓		2	0	0	1	1	1	✓	0	0	1	1	0	0
com.cube.gdpc.isr	health & fitness		✓	0	0	0	1	1	1	✓	0	0	0	1	1	1
com.cyworld.camera	photography			0	0	0	1	1	1		0	0	0	1	1	1
com.devuni.flashlight	tools	✓		0	0	0	1	1	1		0	0	0	1	1	1
com.digisoft.TransparentScreen	entertainment	✓		2	0	0	1	1	1		2	0	0	1	1	1
com.domobile.aplock	tools	✓		0	0	0	1	1	1		0	0	0	1	1	1
com.dropbox.android	productivity	✓		0	0	0	1	1	1	✓	0	0	0	1	1	1
com.ea.game.fifa14_row	games/sports			0	0	0	1	1	1	✓	0	0	0	1	1	1
com.ebay.mobile	shopping			0	0	0	1	1	1		0	0	0	1	1	1
com.facebook.katana	social	✓		0	0	0	1	1	1		0	0	0	1	1	1
com.facebook.orca	communication			0	0	0	1	1	1		0	0	0	1	1	1
com.g6677.andriod.cbaby	games/casual			1	0	0	1	1	1	✓	0	0	1	1	0	0
com.g6677.andriod.chospital	games/casual		✓	1	0	0	1	1	1	✓	0	0	1	1	0	0
com.g6677.andriod.design	games/casual		✓	1	0	0	1	1	1	✓	0	0	1	1	0	0
com.g6677.andriod.pnailspsa	games/casual		✓	1	0	0	1	1	1	✓	0	0	1	1	0	0
com.g6677.andriod.princesshs	games/casual			1	0	0	1	1	1	✓	0	0	1	1	0	0
com.gameclassic.towerblock	games/puzzles	✓		0	0	0	1	1	1	✓	0	0	0	1	1	1
com.gameloft.android.ANMP.GloftDMHM	games/casual			0	0	0	1	1	1		0	0	0	1	1	1
com.game.fruitlegendsaga	games/puzzles			0	0	0	1	1	1		0	0	0	1	1	1
com.gau.go.launcherex	personalization			0	0	0	1	1	1		0	0	0	1	1	1
com.glu.deerhunt2	games/arcade		✓	0	0	0	1	1	1	✓	0	0	0	1	1	1
com.goldtouch.mako	news	✓		1	0	0	1	1	1		1	0	0	1	1	1
com.goldtouch.yynet	news	✓		0	0	0	1	1	1		0	0	0	1	1	1
com.google.android.apps.docs	productivity			0	0	0	1	1	1		0	0	0	1	1	1
com.google.android.apps.translate	tools			0	0	0	1	1	1		0	0	0	1	1	1
com.google.android.youtube	media & video			0	0	0	1	1	1		0	0	0	1	1	1
com.google.earth	travel & local		✓	0	0	0	1	1	1	✓	0	0	0	1	1	1
com.halfbrick.fruitninjafree	games/arcade			0	0	0	1	1	1		0	0	0	1	1	1
com.halfbrick.jetpackjoyride	games/arcade	✓		0	0	0	1	1	1		0	0	0	1	1	1
com.icloudzone.AsphaltMoto2	games/racing			0	0	0	1	1	1	✓	0	0	0	1	1	1
com.ideomobile.hapoolim	finance			0	0	0	1	1	1		0	0	0	1	1	1
com.imangi.templerun2	games/arcade		✓	0	0	0	1	1	1	✓	0	0	0	1	1	1
com.kiloo.subwaysurf	games/arcade		✓	0	0	0	1	1	1	✓	0	0	0	1	1	1
com.king.candycrushsaga	games/arcade		✓	0	0	0	1	1	1	✓	0	0	0	1	1	1
com.skiggle.production	social			0	0	0	1	1	1		0	0	0	1	1	1
com.skype.raider	communication			0	0	0	1	1	1		0	0	0	1	1	1
com.UBI.A90.WW	games/arcade			0	0	0	1	1	1		0	0	0	1	1	1
com.viber.voip	communication			0	0	0	1	1	1		0	0	0	1	1	1
com.whatsapp	communication			0	0	0	1	1	1		0	0	0	1	1	1
Total		17		12	27	1	0				22	13	0	9		
Average							0.99	1.00	0.99 (0.98)					1.00	0.84	0.84 (0.74)

Table 5: Detailed summary of the results of the H2 experiment described in Section 5.4