

Gábor György Gulyás*, Gergely Acs, and Claude Castelluccia

Near-Optimal Fingerprinting with Constraints

Abstract: Several recent studies have demonstrated that people show large behavioural uniqueness. This has serious privacy implications as most individuals become increasingly re-identifiable in large datasets or can be tracked, while they are browsing the web, using only a couple of their attributes, called as their fingerprints. Often, the success of these attacks depends on explicit constraints on the number of attributes learnable about individuals, i.e., the size of their fingerprints. These constraints can be budget as well as technical constraints imposed by the data holder. For instance, Apple restricts the number of applications that can be called by another application on iOS in order to mitigate the potential privacy threats of leaking the list of installed applications on a device. In this work, we address the problem of identifying the attributes (e.g., smartphone applications) that can serve as a fingerprint of users given constraints on the size of the fingerprint. We give the best fingerprinting algorithms in general, and evaluate their effectiveness on several real-world datasets. Our results show that current privacy guards limiting the number of attributes that can be queried about individuals is insufficient to mitigate their potential privacy risks in many practical cases.

DOI 10.1515/popets-2016-0051

Received 2016-02-29; revised 2016-06-02; accepted 2016-06-02.

1 Introduction

People exhibit large behavioural uniqueness. The way how they move around [10], purchase goods [11], configure their browser [12], or browse the web [24] make them unique in very large populations. Even various combinations of seemingly innocuous behavioural attributes, such as in which shops they used to buy goods or where they live and work make them unique [27]. For example, it has been demonstrated that only four spatio-temporal positions are enough to uniquely identify an individual

95% of the times in a dataset of one and a half million users [10].

As many devices record people's behavioural patterns today, it becomes relatively easy for third parties to access such personal information. This has detrimental effect on individuals' privacy, as behavioural attributes can be easily used for the purpose of re-identification and tracking in different contexts. For example, advertiser companies often use behavioural attribute values (e.g., the list of installed fonts [12] or applications [19]) to fingerprint website visitors in order to show them more personalized ads. That is, each time a user visits a website where a 3rd party tracker is embedded, the tracker retrieves the user's fingerprint (such as certain fonts installed on the device [6]) and *link* all his visited sites together [3]. As visited sites can belong to very different topics, the tracker can build accurate interest profiles of visitors in order to show them more personalized ads. However, the list of visited sites can also potentially uncover sensitive information such as sexual behaviour or religion among others. Indeed, most popular website categories that employ such user tracking were found to be porn (15%) and dating (12.5%) [7].

Several countermeasures have been introduced against fingerprinting in the context of web-tracking [7]. A simple solution is to completely disable all techniques which allow a tracker to fingerprint visitors, such as the execution of Javascript and Flash programs in a browser. However, this causes significant degradation of user experience and hence it is a less appealing option. Another, more flexible approach is to enable but *constrain* the usage of such techniques. For example, only *limited number* of system fonts can be used/queried by an entire website [20], or only *limited number* of applications can be opened from another application on a smartphone [2]. Hence, the goal of an adversary is to identify at most a given number of attributes (e.g., fonts or applications), whose values make either a *targeted* user or *every* user as unique as possible in the dataset, i.e., the number of users sharing the same values of these attributes is minimized.

In this paper, we show that the practice of limiting the number of attribute queries per user to mitigate the privacy risks of fingerprinting is often inadequate. We construct efficient greedy fingerprinting algorithms given constraints on the maximum size of the finger-

*Corresponding Author: Gábor György Gulyás: IN-

RIA, E-mail: gabor.gulyas@inria.fr

Gergely Acs: INRIA, E-mail: gergely.acs@inria.fr

Claude Castelluccia: INRIA, E-mail:

claude.castelluccia@inria.fr

User #	A_1	A_2	A_3	A_4
U_1	1	0	1	1
U_2	1	1	1	1
U_3	0	1	0	1
U_4	1	0	1	0
U_5	1	1	1	0
U_6	1	1	0	0

Table 1. Example dataset where each row is a record of a user with binary attributes A_1, A_2, \dots, A_4 .

print. We provide essentially the best fingerprinting algorithms under standard computational assumptions. In particular, we address two general problems. Suppose a positive integer s . First, given a targeted user and all his attribute values, we want to know the subset of attributes with cardinality at most s , called the individual fingerprint of the user, whose values make this user as unique as possible in a given dataset. This is referred to as the *targeted fingerprinting problem* and analogous to the re-identification (de-anonymization) problem when, for example, a user of the Tor Browser is re-identified using some of his/her attributes which can be queried by visited websites. Second, we want to know the subset of attributes with cardinality at most s , called general fingerprint, whose values make each member of a group as unique as possible within the group. This is referred to as the *general fingerprinting problem* and analogous to the linking problem when, for example, the browsing activities of users are tracked, perhaps without re-identifying them.

For example, in Table 1, the individual fingerprint of U_1 with maximum size 2 is $\{A_2, A_4\}$, and the fingerprint value is $\{A_2 : 0, A_4 : 1\}$ which make U_1 unique among 5 other users. However, U_2 has many possible fingerprints with maximum size 2 (e.g., $\{A_2, A_4\}$), but none of them provides a unique fingerprint value to U_2 , because there is always at least another user which remains indistinguishable from U_2 with respect to the fingerprint. Table 2 shows the anonymity sets of all possible subsets of attributes with cardinality 3 (in one anonymity set, users share the same values of these attributes). We find that the only general fingerprint of all the six users, with maximum size 3, is $\{A_2, A_3, A_4\}$, because this is the combination of attributes (with size at most 3) which yields the smallest anonymity sets.

We believe that the solutions of the above fingerprinting problems also have great importance in practice due to the upcoming European General Data Protection Regulation which mandates privacy risk assessment [1]. This requires to measure the re-identification risk

given certain “reasonable” constraints on the adversarial background knowledge. Moreover, the problem of user fingerprinting indirectly also appears in the Article 29 Working Party Opinion 05/2014 on data anonymization techniques [5] which describes the linkability of records concerning the same data subject in an anonymized dataset as a privacy weakness.

Finally, we must note that our problems are different from the one in [21], which is about finding a minimal subset of attributes that provides a certain separation of all the users in a dataset. That is, given a separation value α , the minimum subset of attributes which provides separation α (i.e., it separates at least α fraction of all possible record pairs) is approximated in [21]. By contrast, in our problems, we are given a threshold s and we aim to give a subset of attributes with size at most s which provides the best separation. Obviously, the larger s the better separation is provided. Moreover, as we describe in Section 4.2, our general fingerprinting algorithm solves some problems in [21] more efficiently than prior algorithms. We also note that the naive approach of using a limited number of *random* attributes as a fingerprint, which is also employed by prior re-identification studies [4] [10] [11], is inferior to our solution. In particular, these naive approaches build individual fingerprints from only those attributes whose values are set in a record. For example, in Table 1, $\{A_2, A_4\}$ is not a possible fingerprint of U_1 with such a naive approach, as A_2 is not set in the first record. In fact, U_1 has no unique fingerprint with size 2 out of the attributes whose values are 1 in U_1 's record. Even more, in our work, missing attributes such as $\{A_1 = 0\}$ could be a valid individual fingerprint of U_3 .

Our main contributions are summarized as follows:

- We provide essentially the best fingerprinting algorithms when constraints of the size of the fingerprint are provided. In particular, we first prove that both targeted and general fingerprinting are **NP**-hard in general and the best possible polynomial time approximations are greedy heuristics.
- We provide real-life applications where these fingerprinting methods can be used in practice. In particular, we analyze the privacy protection used by many popular services, such as iOS or Tor Browser Bundle, which are based on limiting the number items (e.g., installed apps or fonts) that can be queried on a device for the purpose of fingerprinting. Our aim is to raise awareness about the weakness of such privacy guards.

Subset of Attributes	Anonymity sets
$\{A_1, A_2, A_3\}$	$\{U_1, U_4\}, \{U_2, U_5\}, \{U_3\}, \{U_6\}$
$\{A_1, A_2, A_4\}$	$\{U_1\}, \{U_2\}, \{U_3\}, \{U_4\}, \{U_5, U_6\}$
$\{A_1, A_3, A_4\}$	$\{U_1, U_2\}, \{U_3\}, \{U_4, U_5\}, \{U_6\}$
$\{A_2, A_3, A_4\}$	$\{U_1\}, \{U_2\}, \{U_3\}, \{U_4\}, \{U_5\}, \{U_6\}$

Table 2. The anonymity sets of all possible subsets of attributes (from Table 1) with size 3. The single general fingerprint is $\{A_2, A_3, A_4\}$ whose values make every user unique.

- We evaluate our fingerprinting algorithms on different large real-life datasets. Specifically, we analyze the separation of different individual as well as general fingerprints on several datasets, where a record can contain the list of installed fonts, smartphone applications, or visited locations by an individual.

The code base related to this paper is available at https://github.com/gaborgulyas/constrained_fingerprinting.

2 Model

Suppose that each individual in a population P has a set of items from a larger universe \mathcal{U} with size n . The universe can contain any observable items of individuals (e.g., set of all possible applications installed on individuals' smartphones, set of possible visited locations, etc.). Therefore, each user is represented by a binary vector u , called the user's profile, where $u[i] = 1$ only if the user has the corresponding item $i \in \mathcal{U}$.

Our goal is twofold. First, we want to identify *at most* s items, called individual fingerprint, which separate a single individual from others the most, i.e., whose values provide the largest uniqueness to the individual in the population. We refer to this problem as *targeted fingerprinting*. Second, we want to identify *at most* s items, called general fingerprint, which simultaneously separate all users from each other in a group the most, i.e., whose values make each member of the group as unique as possible within the group. We call this problem as *general fingerprinting*.

In particular, we are given the following constraints.

- (1) We have only access to a subset B of the whole population P , where one can think of B as a subsample of P .
- (2) We can query at most s items of a user to check whether the user has these items in his/her profile or not. Then, our exact problems to be solved are as follows.

Targeted fingerprinting: Given a targeted profile u and a subset B of population P . Identify at most s items $\{i_1, i_2, \dots, i_s\} \subset \mathcal{U}$ such that the number of profiles in B which share identical values at positions i_1, i_2, \dots, i_s with u is minimized.

General fingerprinting: Given a subset B of population P . Identify at most s items $\{i_1, i_2, \dots, i_s\} \subset \mathcal{U}$ such that the number of profiles in B which share identical values at positions i_1, i_2, \dots, i_s is minimized.

Simply put, we intend to compute the uniqueness of different itemsets in B and generalize these results to the whole population. The underlying assumption of this approach is that B is a subset of the population P , and hence if B has sufficient number of users, the uniqueness of different itemsets in B and P are also similar.

3 Hardness of fingerprinting

We show that identifying the fingerprint items for both problems described in Section 2 is **NP-hard**. In both cases, we reduce the problem to the Maximum Coverage Problem which is **NP-hard** [9].

Maximum Coverage Problem: Given a positive integer s and a collection of sets $C = C_1, C_2, \dots, C_m$. Find a subset $C' \subseteq C$ of sets such that $|C'| \leq s$ and the number of covered elements $\left| \bigcup_{C_j \in C'} C_j \right|$ is maximized.

The reductions are as follows.

Targeted fingerprinting: Let u denote the profile of the targeted user. Assign a set of profiles $D_i \subseteq B$ to each item $i \in \mathcal{U}$, where item i of each profile in D_i has value $u[i]$. Let $\mathbb{D} = \{D_i | i \in \mathcal{U}\}$, where $|\mathbb{D}| = |\mathcal{U}|$. The problem is to find at most s sets in \mathbb{D} such that the size of their intersection is minimized. This is identical to finding at most s sets such that the size of the complement of their

intersection is maximized. Following from De Morgan's laws, this is further equivalent to finding at most s sets in $\mathbb{D} = \{\overline{D_i} | i \in \mathcal{U}\}$ such that the size of their *union* is maximized, where $\overline{D_i}$ contains all profiles whose value at i is *not* $u[i]$. This problem is identical to the Maximum Coverage Problem, and hence Problem 1 is **NP-hard**.

General fingerprinting: Let V denote the number of all pairs of profiles from B , where $|V| = \binom{|B|}{2}$. For each item $i \in \mathcal{U}$, let $T_i \subset A$ denote all pairs of profiles which have *different* values at i (i.e., they can be separated by item i). Clearly, $\bigcup_{i \in \mathcal{U}} T_i = V$. Let $\mathbb{T} = \{T_i | i \in \mathcal{U}\}$, where $|\mathbb{T}| = |\mathcal{U}|$. Hence, Problem 2 is equivalent to finding at most s sets in \mathbb{T} such that the size of their union (i.e., the total number of separated pairs of profiles) is maximized, which is again the Maximum Coverage Problem.

4 Algorithms

As the underlying problems for both fingerprinting problems are **NP-hard**, we cannot hope for finding the optimal solutions. However, a simple greedy heuristics, which approximates the solution of the Maximum Coverage problem, provides essentially the best possible approximations for our fingerprinting problems.

The greedy algorithm for the Maximum Coverage Problem iteratively picks the set which covers the largest number of uncovered elements. It can be shown that this algorithm achieves an approximation ratio of $(1 - 1/e) \approx 0.632$ [22]. That is, if OPT denotes the optimal solution (i.e., the maximum number of covered elements by at most s sets), then the greedy algorithm will cover at least $(1 - 1/e) \cdot OPT$ elements. Moreover, it has also been proven that, unless **P = NP**, there is no $1 - 1/e - o(1)$ approximation for the Maximum Coverage Problem [13], which means that the greedy approach is essentially the best-possible polynomial time approximation for this problem.

4.1 Targeted fingerprinting

The first problem is straightforward to approximate with the greedy approach which is described in Algorithm 1. Instead of maximizing the number of users which disagree on the values of fingerprint queries with the targeted user u , we minimize the number of users which *agree* on those values with u .

We iteratively squeeze the anonymity set of u , denoted by $anon_set$ in Algorithm 1, which is the set of users sharing identical values at all items of the fingerprint queries K with u . Initially, $anon_set$ is composed of all the users in B . Then, in each iteration, we select the item i_{sep} which separates u from the other profiles in the anonymity set the most, i.e., minimizes the number of profiles which agree at all items of K with u . For this purpose, we pre-compute a table, denoted as $users$ in Line 2, which maps each item $i \in \mathcal{U}$ to the list of profiles which has the value of $u[i]$ at position i . The computation of this table takes time of $O(|B||\mathcal{U}|)$. Then, we can easily identify the most separating item in each iteration by the intersection of the anonymity set of u and each list of table $users$. Fast implementations of the intersection of sorted integers are described in [18]. If users in the anonymity set are no further distinguishable by any item, we terminate and return K as the final fingerprint (in Line 7). Otherwise, we continue the separation as long as $|K| \leq s$.

This solution has a storage complexity of $O(|B||\mathcal{U}|)$ and computational complexity of $O(s|B||\mathcal{U}|)$.

Algorithm 1: Greedy Approximation for Targeted Fingerprinting

Input: Fingerprint size s , Universe of items \mathcal{U} , Dataset B , Targeted user's profile u
Output: Fingerprint items K

```

1  $K := \emptyset$ 
2  $users[i] := \{u' | u' \in B \wedge u[i] = u'[i]\}$  for all  $i \in \mathcal{U}$ 
3  $anon\_set := B$ 
4 repeat
5     // if no further separation is possible
6     if  $\min_{i \in \mathcal{U} \setminus K} \{|users[i] \cap anon\_set|\} = |anon\_set|$ 
7         then
8             break
9      $i_{sep} := \arg \min_{i \in \mathcal{U} \setminus K} \{|users[i] \cap anon\_set|\}$ 
10     $anon\_set := anon\_set \cap users[i_{sep}]$ 
11     $K := K \cup \{i_{sep}\}$ 
12 until  $|K| = s$  or  $|anon\_set| = 1$ 
13 return  $K$ 

```

Example 1. Suppose that $u = U_5$ in Table 1 and $s = 3$. First, we select the attribute which makes U_5 the most unique. U_5 shares the same values of A_1, A_2, A_3 and A_4 with 4, 3, 3, and 2 other users, respectively, hence we select A_4 in the first iteration. As a result, U_4 and U_6 still remain indistinguishable from U_5 . In the next iteration, we select the attribute which distinguishes U_5 from U_4 and U_6 the most. We have two options: A_2 or

A_3 , as A_1 are identical for all of them. We select, say, A_2 , which means that U_5 and U_6 are still indistinguishable. Hence, in the final iteration, we select A_2 , which yields $\{A_2, A_3, A_4\}$ as the individual fingerprint of U_5 .

4.2 General fingerprinting

The naive greedy algorithm, which includes enumerating and storing each pair of profiles distinguished by each attribute, has a storage complexity of $O(|\mathcal{U}||B|^2)$ and computational complexity of $O(s|\mathcal{U}||B|^2)$. Instead, we design a greedy algorithm whose storage complexity is $O(|\mathcal{U}||B|)$ and its computational complexity is $O(s|\mathcal{U}||B|)$.

Our proposal is shown in Algorithm 2. We iteratively partition the set of profiles in B into anonymity sets, where all profiles in an anonymity set have identical values at all items of K (i.e., they are not distinguished by any items in K). Our goal is to compute a partitioning where each partition (i.e., anonymity set) is as small as possible.

To do so, in each iteration, we compute a separation metric of each item for the current partitioning C , which is the total number of separated pairs of profiles over all partitions (in Line 7-11). For example, if a partition is split into sub-partitions P_1 and P_2 by item i (where all profiles in a partition share identical value at position i), then the separation of i with respect to P is $|P_1| \times |P_2|$. The total separation of i is the sum of separations over all partitions. Then, the item with the largest separation, denoted by i_{sep} , is selected (Line 12). Finally, each partition in C is split by i_{sep} into two sub-partitions (Line 17-21) such that a sub-partition contains all profiles from C which have identical values at position i_{sep} . To speed up the computation of the separation metric, we pre-compute a table called *items* which maps each profile u to the set of items which are contained by u (in Line 2 of Algorithm 2). Again, this pre-computation runs in $O(|\mathcal{U}||B|)$.

Example 2. Suppose that $s = 3$ and we want to compute the general fingerprint of Table 1. The operation of the algorithm can also be illustrated by a binary tree which is shown in Figure 1. In the first iteration, we compute the separation metric of all attributes A_1, A_2, A_3, A_4 which are $5 \times 1, 4 \times 2, 4 \times 2,$ and $3 \times 3,$ respectively. Hence, we select A_4 which separates all users into two partitions $\{U_1, U_2, U_3\}$ and $\{U_4, U_5, U_6\}$. In the second iteration, we again compute the separation metric of A_1, A_2, A_3 over the two partitions, which are $2 \times 1 + 0 = 2,$

$2 \times 1 + 2 \times 1 = 4, 2 \times 1 + 2 \times 1 = 4,$ respectively. Hence, we can select A_2 or A_3 . We select, say, A_2 , which produces partitions $\{U_1\}, \{U_2, U_3\}, \{U_4\}, \{U_5, U_6\}$. Next, we compute again the separation metrics of A_1 and A_3 over the new partitions which are $0 + 1 \times 1 + 0 + 0 = 1$ and $0 + 1 \times 1 + 0 + 1 \times 1 = 2,$ respectively. Therefore, we finally select A_3 . As 3 attributes have been selected, we terminate and return $K = \{A_2, A_3, A_4\}$ as the general fingerprint.

Notice that the maximum number of partitions in any iterations is at most $|B|$ (i.e., the number of all records), and we only maintain the partitioning resulted by the last iteration. The complexity of computing the separation metric of each attribute takes $O(|\mathcal{U}||B|)$ (Line 7-11) in each iteration, as set-memberships can be checked in $O(1)$ using hash-maps. Similarly, the separation itself (Line 17-21) runs in $O(|\mathcal{U}||B|)$, which means that the total complexity is $O(s|\mathcal{U}||B|)$ (Line 5-22).

Algorithm 2: Greedy Approximation for General Fingerprinting

Input: Fingerprint size s , Universe \mathcal{U} , Dataset B
Output: Fingerprint items K

```

1  $K := \emptyset$ 
2  $items[u] := \{i \mid i \in \mathcal{U} \wedge u_i = \text{True}\}$  for all  $u \in B$ 
3  $users[i] := \{u \mid u \in B \wedge u_i = \text{True}\}$  for all  $i \in \mathcal{U}$ 
4  $C := B$  // anonymity sets
5 while  $|K| < s$  do
6   // compute the separation of all items
7    $separation[i] := 0$  for all  $i \in \mathcal{U}$ 
8   for  $S \in C$  do
9     for  $i \in \bigcup_{u \in S} items[u]$  do
10       $t[i] := |\{u \mid u \in S \wedge i \in items[u]\}|$ 
11       $separation[i] := separation[i] + t[i] \cdot (|S| - t[i])$ 
12    $i_{sep} := \arg \max_{i \in \mathcal{U} \setminus K} separation[i]$ 
13   // if no further separation is possible
14   if  $separation[i_{sep}] = 0$  then
15     break
16   // splitting all partitions with  $i_{sep}$ 
17    $T = \emptyset$ 
18   for  $S \in C$  do
19      $V := \{u \mid u \in S \wedge i_{sep} \in items[u]\}$ 
20      $T := T \cup \{V, \{S \setminus V\}\}$ 
21    $C := T$ 
22    $K := K \cup \{i_{sep}\}$ 
23 return  $K$ 

```

Finally, we note that Algorithm 2 can also be adapted to compute the minimum number of items which separate all users in B (this is referred to as the

Minimum Key Problem in [21]). To do so, we only need to set s to $|\mathcal{U}|$ in Algorithm 2, which means that the running complexity becomes $O(|\mathcal{U}|^2|B|)$. By contrast, the best prior greedy solution of this problem, which is also proposed in [21], runs in $O(|\mathcal{U}|^3|B|)$. However, this is impractical if the number of attributes/items is large. Although several probabilistic relaxations of this problem is considered in [21] in order to reduce running complexity, they also degrade the accuracy of the solution.

5 Applications

In this section, we describe some potential applications of targeted and general fingerprinting in practice. The first application relates to fingerprinting smartphone applications, which has also been partially studied in [4]. The second application shows how to fingerprint users using the list of (non-)installed fonts in order to track them while they are browsing the web. Finally, we show a use-case of mass de-anonymization in location datasets using our fingerprinting techniques, which can serve as a potential investigation tool for authorities.

5.1 Smartphone app fingerprinting

It has been shown that knowing the presence of four randomly chosen applications on a device are enough to re-identify users in a dataset of 54,893 users [4]. Moreover, installed apps also reveal a lot about the owner's interests, and Twitter also used this information for targeted interest-based advertising among others [19]. In this section, we show how a third party can use such information for tracking. The main idea is to install a malicious application on the user's device which executes the fingerprint queries and identifies the fingerprint value of the user. As some operating systems, such as iOS, limit the number of apps whose existence can be checked by another app, the queries should be optimized such that their result distinguish as many users as possible. In particular, we envision a re-identification scenario using targeted fingerprinting (see Fig. 2a), when the browsing activities of some targeted users are tracked, and also a linking scenario using general fingerprinting (see Fig. 2b), when the browsing activities of multiple users are tracked without re-identifying any of them.

In both cases, when the (individual or general) fingerprint apps are provided to the malicious app, it

checks their existence on the device (i.e., constructs the fingerprint value). Then, it opens a legitimate URL in the browser which redirects to a third party tracker in order to generate and store a tracking cookie of the tracker as a visited first party in the browser¹. This cookie contains the user/device identifier which is the fingerprint value. When the same browser visits a site having this tracker, the user will be tracked using the previously stored cookie (and also re-identified in case of targeted fingerprinting).

We must note that, in contrast to traditional web cookies, the identifier stored in the cookie is strictly bound to a persistent, behavioural attribute (i.e., the list of installed/non-installed apps). This means that even if the user deletes the cookie, he can still be re-linked to the same identifier by re-generating the same cookie based on possibly the same list of (non-)installed apps. The identifier even persists, if the application is removed and then re-installed, or other applications are installed that can generate the same identifier.

Next, we provide the details of the attack, which consists of three phases: (1) *Identify the fingerprint apps (queries) of the targeted user (targeted fingerprinting) or a group of users (general fingerprinting)*, (2) *Retrieving the list of fingerprint apps on users' device*, (3) *De-anonymize the targeted user and/or track user activities*. As the third phase has been detailed before, we will focus on the first two phases.

5.1.1 Identifying the fingerprint apps

For targeted fingerprinting, the malicious 3rd party tracker, called as adversary henceforth, first needs to obtain the list of (non-)installed apps of a known, targeted user. This is relatively easy by consulting with some social media where users often post the applications they use, or simply retrieving the app list directly from the user's device (if the user has iOS 9, the adversary may get access to the full list of installed apps before the user upgrades to iOS 9, see below for details). As another work suggests [4], the adversary does not need to get the complete list of applications, but

¹ In order to circumvent the third-party blocking of Safari, the legitimate URL specifies a malicious website which redirects to the third-party tracker. The tracker then becomes a first party and can set a tracking cookie in the browser. Finally, the tracker redirects to the initial website, and the device can be tracked across websites.

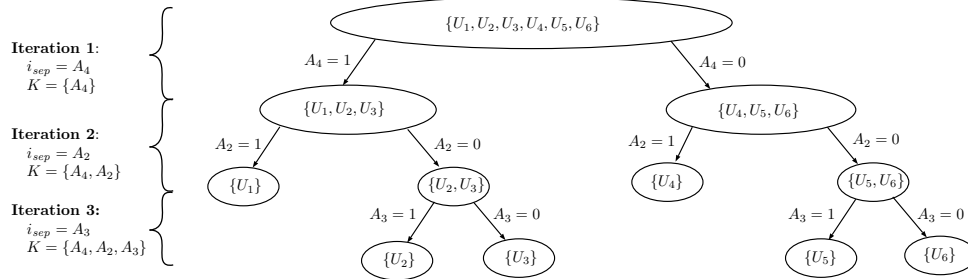


Fig. 1. General fingerprinting on the dataset of Table 1: Greedy algorithm builds a binary tree of partitions, where the leaves constitute the final anonymity sets of the fingerprint $K = \{A_2, A_3, A_4\}$.

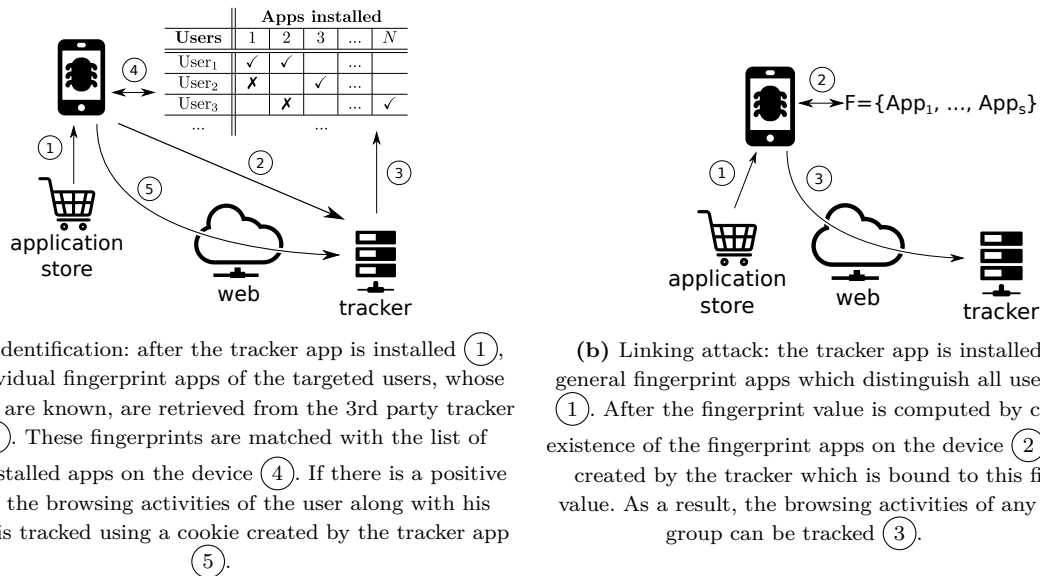


Fig. 2. Targeted fingerprinting (for re-identification) and general fingerprinting (for linking) using the list of installed apps on a smartphone.

only 4 random apps are sufficient to re-identify the user with 95% of the times. We show later that, by using the greedy optimization algorithm, described in Section 4.1, the adversary can re-identify most users with even fewer queries.

For both targeted and general fingerprinting, the adversary needs to have access to a subset of user profiles, denoted by B in Algorithms 1 and 2, in order to identify the fingerprint apps of the users to be tracked. In particular, for targeted fingerprinting, the adversary identifies at most s apps, using Algorithm 1, whose existence need to be checked on the targeted user's device to construct his/her individual fingerprint. As we detail below, the value of s changes depending on which operating system the app is built for. In certain cases, iOS allows to query at most 50 apps, while in other cases the queries need to be pre-defined and cannot be changed

later when the app is installed, which is reminiscent of our general fingerprinting problem.

For general fingerprinting, the dataset is used to identify at most s fingerprint apps, using Algorithm 2, whose existence need to be checked on *every* group member's device to construct a fingerprint for each group member which sufficiently distinguish them. This dataset can contain the profiles of the group themselves, or another subset of the *same* population. After calculating the user/device identifier, this could be used as an identifier for web tracking; we show how in Fig. 3 (with a proof-of-concept demo).

5.1.2 Retrieving the fingerprint apps on iOS

Equipped with the list of the fingerprint apps, the malicious application needs to check their existence on

the device. To the best of our knowledge, Android allows to retrieve the complete list of installed apps, but Apple has introduced some limitations with the release of iOS 9 [2]. On iOS, applications can use the `canOpenURL()` call to determine if another application is available on the system (i.e., whether it can be called). That is, each fingerprint app query corresponds to a call of `canOpenURL()` specifying a URL. The call returns a boolean value indicating whether or not the specified URL's scheme can be handled by some app installed on the device. However, Apple introduced two restrictions that apply when the program is compiled on iOS 9 [2]. Applications built for systems before iOS 9 but executed on iOS 9 can invoke `canOpenURL()` with at most 50 distinct apps. This number is reset each time the user re-installs or upgrades the app. Applications built for iOS 9 and also executed on iOS 9 have to declare in advance which other applications they would like to check with `canOpenURL()`, making abuses noticeable well in advance. Every `canOpenURL()` call which is over the limit of 50, or specifies a non-declared app, will return a negative response.

We emphasize again that these limitations are in effect only when a user builds an app on iOS 9, leaving users of older systems vulnerable. The proportion of such systems was measured to be 23% according to a statistic revealed by the Apple Store (as of February 22, 2016) ².

5.2 Font fingerprinting

Fingerprinting projects, such as the Panopticlick³ [12] and the cross-browser fingerprinting project⁴ [6] have revealed that websites can effectively query browser properties to track their visitors. In particular, it has been shown that detecting fonts is an important technique to track users as their entropy is one of the highest among the available identifiers [12], and they can be queried independently from the browser [6]. Font detection mechanisms were also shown to be used in device tracking mechanisms [23].

² Statistic published by the Apple Store: <https://developer.apple.com/support/app-store/>

³ Website: <https://panopticlick.eff.org>

⁴ Website: <https://pet-portal.eu/fingerprint/>; newer version available at: <https://fingerprint.pet-portal.eu>

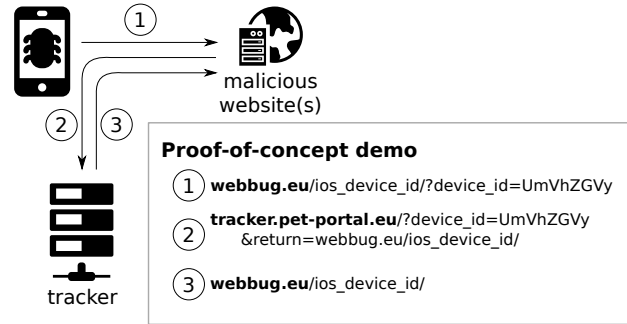


Fig. 3. Example on how the device identifier could be leaked to the web. The malicious application opens a URL that seems to be legitimate, but it contains the hidden or obfuscated device identifier (1). In order to circumvent the potential third-party blocking of Safari, the malicious websites redirects to a third-party tracker (2). The tracker then is now becomes a first party and can set a tracking cookie. Finally, the tracker redirects to the initial website (3), while it also remains able to track the devices across websites.

One cutting edge web browser that promises anonymity is called the Tor Browser Bundle ⁵ (TBB) or the Tor Browser, whose anonymity guarantees have been widely studied [8, 17]. Up to version 5.5, TBB developers tackled font-based tracking by introducing two limits: one on the number of font load trials, and one on the number of successful font loads [20]. These parameters are called `browser.display.max_font_count` and `browser.display.max_font_attempts`, and by default they are set to a value of 10. These options can be changed by accessing `about:config`. Detecting the existence of fonts on a system is relatively easy using Javascript, which can be embedded on the visited website: it can be measured if the rendered text has a different width and height compared to the system default ⁶. Notice that TBB allows the execution of Javascripts by default due to their prevalence on the web today. TBB addresses two types of attacks, among others, which we analyze in this paper: de-anonymization (see Section 3.1.2 in [20]) and linking attacks (see Section 3.1.4 in [20]). Our attack schemes are depicted in Fig. 4.

De-anonymization occurs when an attack correlates TBB and non-TBB activities of a targeted user. First, the adversary collects the available fonts on the user's system when the user visits a website with a regular, non-TBB browser, thus the list of fonts can be associ-

⁵ The TOR project website can be found here: <https://www.torproject.org>

⁶ Example: <http://www.lalit.org/lab/javascript-css-font-detect/>

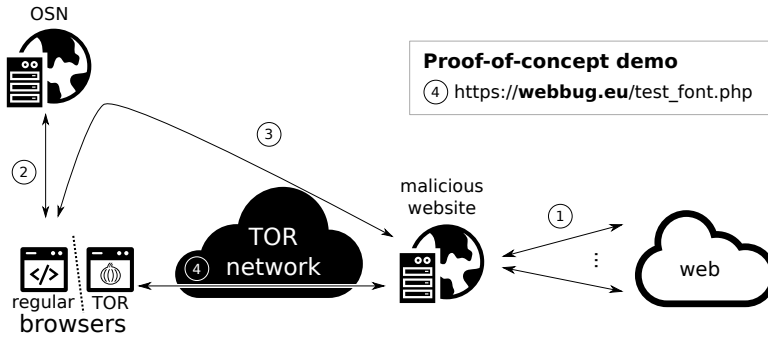


Fig. 4. Attacks against TOR browser. Initially, the malicious websites (or collaborating sites) build a large dataset of (font) profiles by installing trackers (or web bugs) on websites which are also accessed without TOR (1). In case of *targeted fingerprinting*, when the user visits a site containing a malicious tracker (2), the adversary gets informed by its trackers and retrieves the complete profile of the targeted user (3). Then, using the collected dataset of font profiles, the adversary can compute the fingerprint fonts of the user, which can be later queried in order to de-anonymize the user when he visits a site through TOR containing a malicious tracker (4). In case of *general fingerprinting*, the adversary computes the general fingerprint using the collected dataset of font profiles, and queries these fonts from each TOR visitor (4).

ated with the user’s identity. We note that this does not mean that the site is under a total control of the adversary. This can also be a site where a malicious tracker (or web bug) is installed. Then, the adversary selects at most 10 fingerprint fonts whose existence need to be checked on the user’s system to re-identify him/her. When the user visits this or another malicious website via Tor Browser, this site queries these fingerprint fonts through Javascript. As a result, “private” (TBB) and “public” (non-TBB) activities of the web user can be linked. De-anonymization is analogous to our targeted fingerprinting problem, where the maximum number of fingerprint queries is limited to 10.

Linking attack refers to linking the activities of a user together within TBB. As opposed to de-anonymization, the identity of the user is not known (i.e., the attack does not require a visit to a site with a regular browser), and the goal is to track the user using TBB. In this case, the adversary needs to identify at most 10 fingerprint fonts which can distinguish *any* TBB user from the rest. This means that a malicious tracker queries the *same* set of 10 fonts of each TBB user visiting a site, and the results of these queries serve as a fingerprint for each user. This is analogous to our general fingerprinting problem, where the maximum number of fingerprint queries is limited to 10. Notice that while these fingerprints may not be sufficiently unique to every user in the whole TBB community, they should be quite unique per site. According to a measurement based on TBB usage data of 2012-2013 [15], the average number of daily Tor Browser users per country was well below 1%, e.g., 0.08% in Europe, meaning that even visi-

tors of well known websites receiving a large daily traffic could be easily identified even if they would use TBB.

In order to identify the fingerprint fonts, the adversary needs to have access to a subset of all profiles in the TBB population, where a profile contains the list of installed and non-installed fonts of a user.

Note on deprecation: In the beginning of our investigation, we found out that the font limitation was not working at all in the Tor Browser, making TBB users vulnerable to the above attacks without limitation (the current stable version was TBB 5.0.4 at that time). We published a blog post to demonstrate the feasibility of the above attacks [16], and informed the developers about this issue. It turned out that developers had already been working on other workarounds to defeat font fingerprinting [14], and the lack of the above limitation on the number of font queries was due to implementation difficulties in newer Firefox versions. The new countermeasure, which is shipping the Tor Browser with a pre-defined set of fonts that the browser can only use, was introduced two months later with TBB 5.5 at January 26, 2016. While this makes our fingerprinting techniques deprecated in the context of TBB, we believe that our work can be a useful warning for developers who may use similar countermeasures in the future, as the case of application fingerprinting also shows.

5.3 Mass de-anonymization in location datasets

We can also use fingerprinting techniques in a quite different context than online tracking, that is, to de-anonymize a location database (e.g., a Call Detail Record dataset). Suppose that the adversary has access to a weakly anonymized location dataset containing the set of visited locations per individual in a city (time information is omitted for the sake of simplicity). For example, this dataset contains pseudonymized records (i.e., only direct personal identifiers are removed), but location information are kept intact. The goal of the adversary is to re-identify as many users as possible in this dataset, i.e., to perform a *mass de-anonymization attack without knowing any locations of the victims a priori*⁷, but only their personal photos (e.g., by crawling a social media of the city). The plan of the adversary is to mount at most s cameras at different spots of the city, and re-identify all individuals observed by any of these cameras by using some face recognition mechanism and the photos that she has. Knowing the fact that an individual does (not) visit these camera spots (i.e., the visiting pattern of individuals with respect to only these spots), the adversary locates the record which has the observed visiting pattern in the anonymized dataset, and learn all other unknown visited places of the individual.⁸

In order to maximize the number of de-anonymized users in the dataset, the adversary wants to identify at most s camera spots which can distinguish *all users* in the anonymized dataset as much as possible. In particular, the values of *any* record at these spots should be unique in the dataset, or at least the number of records sharing the same values at these spots should be minimized. This goal is reminiscent of our *general fingerprinting problem*. Hence, the adversary runs our greedy heuristic in Algorithm 2, with input s and the anonymized dataset, in order to identify these fingerprint spots.

Although deploying cameras at any spot of a city is unlawful and easily detectable, we believe that the above scenario is not far-fetched at all. Authorities, who have the right to install or access such surveillance cameras, may be interested in the above “privacy-friendly” approach of learning the trajectories of suspects. Indeed,

authorities do not need to ask telecom companies to reveal the exact identity of any of their customers, but only their anonymized CDRs.

6 Experimental results

In this section, we compute the uniqueness of targeted and general fingerprints on real-world datasets. In particular, we consider three real-world datasets, each corresponding to an application described in Section 5, which can be used by the adversary to identify the fingerprint items of users. We report the size of the anonymity sets of all targeted and general fingerprints for each dataset.

6.1 Datasets and ethical considerations

6.1.1 Smartphone applications

6.1.1.1 Dataset

The dataset comes from the Carat research project [25]. The dataset includes data from 54,893 Carat Android users between 11/03/2013 and 15/10/2013 [26]. During this period, the Carat app⁹ was collecting the list of running apps (and not the list of all installed apps) on users’ devices when the battery level changes. As collecting the list of running apps multiple times over more than 7 months is likely to sum up to the set of all installed apps of a user, we consider a record as the set of installed applications in this paper, even if a record might not be the complete set of installed apps all the time.

We removed system apps from all records because they are common to all users. Without system apps, our analyzed dataset contains 92,210 different applications whereas the total number of apps available on the GooglePlay were around 1 million during this time¹⁰. Table 3 summarizes the main characteristics of our dataset.

6.1.1.2 Ethical considerations

This dataset contains the list of installed smartphone applications per user. The data were collected with the users’ consent, and they were explicitly informed that

⁷ Otherwise, the adversary has good chance to locate the victim’s record based on [10]

⁸ It is not reasonable to assume that the adversary can mount any number of cameras due to budget constraints.

⁹ <http://carat.cs.helsinki.fi>

¹⁰ http://en.wikipedia.org/wiki/Google_Play

# of records	54,893
# of all apps in the dataset	92,210
Maximum record size	541
Minimum record size	1
Average record size	42
Std.dev of record size	39

Table 3. Characteristics of the smartphone app dataset

# of records	43,656
# of all fonts in the dataset	81,195
Maximum record size	6,072
Minimum record size	1
Average record size	143
Std.dev of record size	118

Table 4. Characteristics of the font fingerprint dataset

their data could be used and shared for various research projects. In fact, the Carat privacy policy (available at <http://carat.cs.helsinki.fi>) clearly specifies that “Carat is a research project, so we reserve the right to publish our results online and in academic publications. We also reserve the right to release the data sets into the public domain.” Also, the dataset was shared with us by the Carat team in a pseudo-anonymised form. In particular, identifiers were removed, and each application name was replaced with its SHA1 hash. It contained 54,893 records [26], i.e. one record per user. Each record is composed of the list of applications installed and run by the user. Furthermore, the data sharing agreement that we signed, stipulated that we cannot use the dataset to deanonymize the users in the dataset.

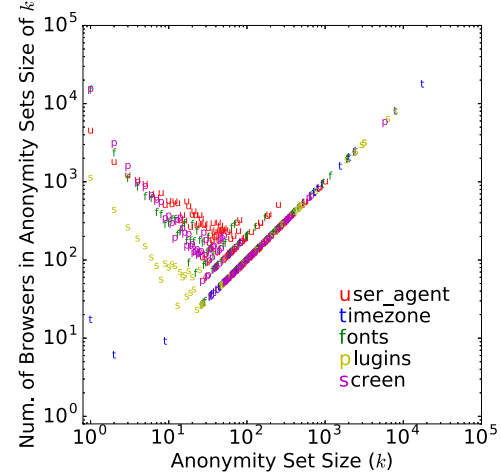
6.1.2 Fonts used by web browsers

6.1.2.1 Dataset

This dataset contains the list of installed fonts, the screen size and available screen size per user. The data collection was done on fingerprint.pet-portal.eu where users could test and compare their OS fingerprint in different browsers. Data was collected with the purpose of analyzing OS fingerprint uniqueness. In particular, users were advised to run the test in different browsers which collected the list of installed fonts, among others, and compare the calculated fingerprint if they match. The website have collected more than 107,102 font profiles between 05/04/2012 and 22/09/2015, which do not necessarily belong to distinct users due to the purpose of the project. In our experiments, we only used a sub-

	Panoptick	current
User agent string	10.0	7.18
Timezone	3.04	2.23
All fonts	13.9	7.79
Plugins	15.4	7.91
Screen	4.83	3.34

(a) Comparison of entropy values for each attribute.



(b) Number of users w.r.t. the sizes of the anonymity set they belong to.

Fig. 5. Descriptive attributes of our datasets that also enable the comparison with previous work in [12].

set of users and attributes that were collected, and the data was cleaned as follows.

In order to remove possible duplicates (a user could visit the fingerprinting site multiple times), several records were removed by comparing hashed IPs and stored usernames. If username was not given by a person who conducted the test, the site stored the font profile of the user in an evercookie¹¹. Given a list of profiles having the same IP address where creation dates were closer than 10, only the first entry was preserved. All record duplicates having the same username were removed regardless of date differences.

The resulted database had 43,656 records, which contained the list of detected fonts, and the screen and available screen sizes¹² of each user. As this is the first

¹¹ Evercookies: <http://samy.pl/evercookie/>

¹² Referring to JavaScript properties `screen.width` & `screen.height` and `screen.availWidth` & `screen.availHeight`

paper using the given dataset, we provide descriptive metrics of the data in order to help understanding our results. We also provide attributes of the original Panoptick dataset for comparison, and Fig. 5b is also easily comparable with Fig. 3 in [12]. This figure confirms that our dataset has been cleaned properly. The main properties of the dataset is described in Table 4.

6.1.2.2 Ethical considerations

The data were collected with users' consent on a website that raises awareness on tracking using OS (or device) fingerprinting. On the website it is clearly described that all collected information is treated confidentially, and the goal of the data collection is to analyze uniqueness of generated OS fingerprints. Types of data that are collected is also clearly communicated, from which we used screen size, available screen size, and available fonts per user. The data we used was shared with us with the operators of the fingerprint website, and our agreement allowed us to conduct the uniqueness study we present in this paper.

6.1.3 Location dataset

6.1.3.1 Dataset

This dataset contains the list of visited locations per user in a large European city. We use a CDR (Call Detail Record) dataset provided by a cell phone operator in Europe. A cell tower is visited by an individual, if the operator has a recorded event at the tower related to the individual over the observed period (01/09/2007 - 17/10/2007). An event can be an incoming/outgoing call or message to/from the individual. The dataset contains the events of 4,427,486 users at 1,303 towers within the administrative region of the city, where the GPS coordinates of all the towers are also available. In our dataset, a record contains *only* the set of towers that are visited by a user over six weeks, i.e., it does *not* contain the time of visits.

The average number of individuals per tower over this period was 38,817 with a standard deviation of 50,911. The total area of the city which is covered by cell towers is 128.1 km². The main characteristics of our dataset and the cell towers are shown in Table 5.

6.1.3.2 Ethical considerations

The dataset was shared with us by the telecom operator in a pseudo-anonymised form. Each record in this

Dataset size	4,427,486
# of all cell towers	1303
Maximum record size	422
Minimum record size	1
Average record size	11.42
Std.dev of record size	17.23
Total area of all cells	128.1 km ²

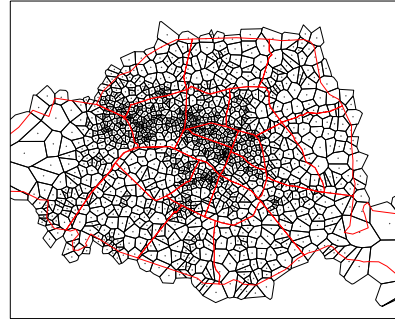


Table 5. Characteristics of our location dataset (up) and Voronoi-tessellation of cell towers (down). Red lines denote the boundaries of districts.

dataset represented an event belonging to a single user with a pseudonymized phone number, i.e., the last 8 digits of the original phone number were changed to a randomly chosen 8 digit number. The record included the GPS location of the cell tower, the type of the event (incoming/outgoing SMS or call) and the timestamp of the event. Records (events) belonging to the same individual contained identical pseudonymized phone number, hence all events of an individual could be linked together. We signed a data sharing agreement with the operator which stipulated that we cannot use the dataset to deanonymize the users in the dataset.

6.2 Targeted fingerprinting

6.2.1 Smartphone users by applications

We run Algorithm 1 on the smartphone application dataset, described in Section 6.1.1, to generate an individual fingerprint for each user in our dataset with different constraints on the maximum size of the fingerprint. When $s = 50$ (this threshold comes from the iOS countermeasure detailed in Section 5.1.2), on average, only 2.3 app queries per user provided the fingerprint value. The exact distribution of fingerprint length is shown in Fig. 6. Our greedy algorithm generated fingerprints shorter than 50 apps for all users. Moreover,

with the exception of 27 users, the generated fingerprints consist of at most 10 apps. The reason is that our algorithm can terminate before the fingerprint size reaches s (see Line 7 in Algorithm 1). We note that a user who has a fingerprint shorter than s is not necessarily unique with this fingerprint in the dataset. For example, a user may have fewer than s apps installed on his phone which are also installed by other users in the dataset.

For $s = 50$, we find 4,098 users (7.47%) in our dataset who cannot be uniquely re-identified and therefore belong to some anonymity sets with size greater than 1. We visualized the sizes of anonymity sets and their population in Fig. 8a. This figure shows that 96.15% of users are almost unique (the size of their anonymity sets is at most 3), while only 3.02% of users fall into anonymity sets with size at least 10. (Note: large anonymity sets concerning only a couple of users are possible as these anonymity sets are not disjoint.)

We also measured the size anonymity of sets when $s = 2$ (instead of 50); results are reported in Fig. 8b. The proportion of almost unique users decreased to 85.54%, eventually showing that the current privacy protection method on iOS is ineffective against fingerprinting users. In addition, this uniqueness value is larger than the one in [4], which was about 75% for 2 apps. The difference, as we have discussed in Section 1, is due to the fact that only existing apps of users are selected to measure their uniqueness in [4], meanwhile we optimize the queries for the targeted user's profile and they can also include the queries of non-existing apps.

Finally, we note that similar uniqueness is expected in even larger datasets based on the same reasoning as in [4]. With a growing number of users, the number of available applications should also increase with a lower pace. This means that the number of possible combinations of apps also increases exponentially with each new application added to the dataset.

6.2.2 TBB users by font detection

We also run Algorithm 1 on the font dataset, described in Section 6.1.2, to measure the effectiveness of fingerprinting using fonts. We believe that, despite that non-Firefox browsers are not excluded from our dataset, our results are generic. While there are minor differences of JavaScript font detection in different browsers, this method is reliable enough to detect the majority of fonts in all browsers. As a result, a record characterizes the OS of a user rather than his browser. We also studied the uniqueness of fingerprints when they are extended with

the screen resolution and available screen size. While the Tor Browser makes various efforts to hinder access to this information, it cannot block full access [20]; for example, it warns users of related privacy hazards in case the window of a Tor Browser is maximized.

The distribution of fingerprint length is shown in Fig. 7. Results are better for user privacy compared to our app dataset, as only 53.8% of users have fingerprints no longer than the font load limit of $s = 10$ imposed by the Tor Browser (see Section 5.2). Peaks appearing for longer length are likely to be due to pre-set OS configurations that were not altered at all, or only slightly. For example, these can be computers running in institutions or at companies, where configurations are very similar to each other. Barely modified fresh installs could also be behind this phenomena.

As font fingerprints are not necessarily unique (due to the same reasons discussed in Section 6.2.1), we report the sizes of anonymity sets for each user. Results are shown in Fig. 9a. Detection of 10 fonts allows unique identification of 13,155 users in our datasets (30.13%), and 39.26% of the users are almost unique (i.e., has anonymity set size at most 3). This roughly means that one third of TBB users at a given service can effectively be de-anonymized.

We also report the size of anonymity sets if screen resolution and available screen size are also added to the fingerprint (see Fig. 9a and 9b). These results show that these extra attribute values significantly increase the uniqueness of users. While we find that the proportion of almost unique users is roughly doubled when available screen size is added to the fingerprint, probably this would not be the case for TBB users. In particular, TBB users are discouraged to install any browser extensions and other tools which usually change the available screen size.

We additionally examine if decreasing the limit s is a remedy for the situation. Our negative findings show that when the font load limit is decreased to 5, significant uniqueness can be still observed (see Fig. 10 for details). For example, detecting solely 5 fonts makes still 30.58% of users unique (Fig. 10a).

6.3 General fingerprinting

We run Algorithm 2 on the three datasets to compute general fingerprints and partition users into anonymity sets, where users have the same fingerprint value in each set. The average size of anonymity sets for small fingerprint sizes (less than 20) are shown in Fig. 11. The aver-

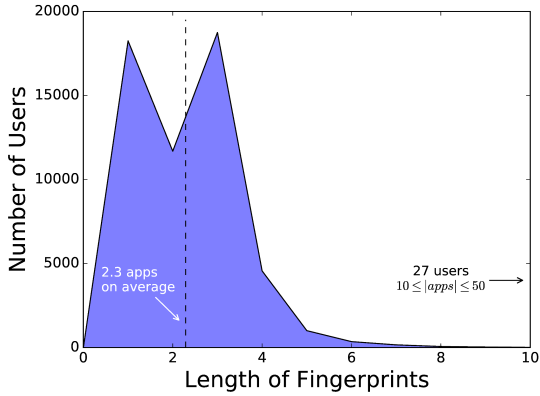


Fig. 6. Distribution of the length of individual fingerprints for the app dataset. Almost all users have fingerprint size smaller than 10, with only 27 exceptions.

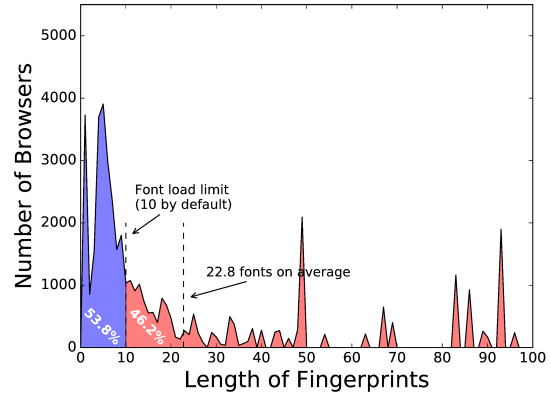
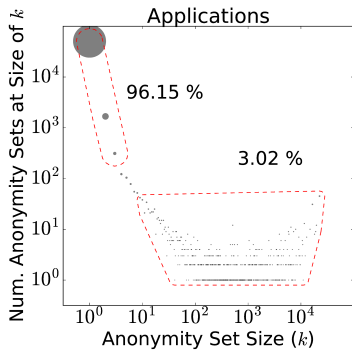
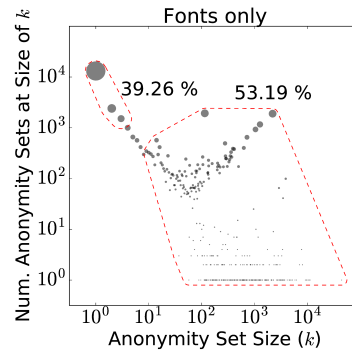


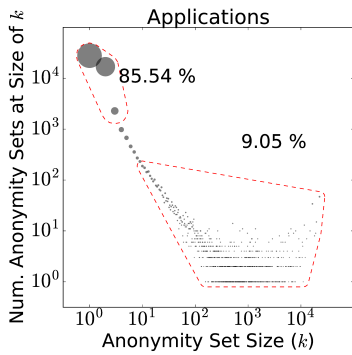
Fig. 7. Distribution of the length of individual fingerprints for the font dataset. 53.8% of users has a fingerprint smaller than the constraint 10.



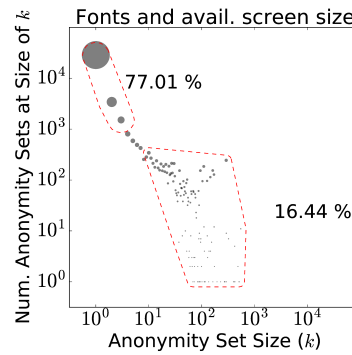
(a) $s = 50$



(a)



(b) $s = 2$

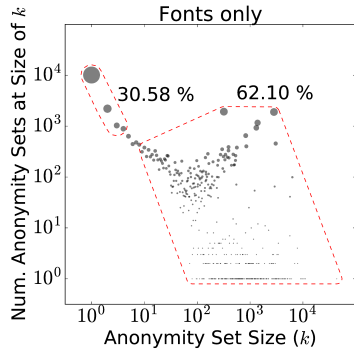


(b)

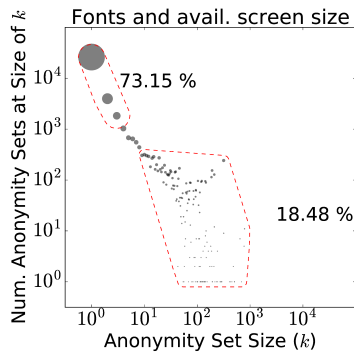
Fig. 8. Targeted fingerprint attacks with the app detection limit defined by Apple (a), and decreased further to only 2 applications (b). (Note: dot sizes is proportional to the proportion of users having a particular anonymity set size.)

Fig. 9. Comparison of different targeted fingerprinting methods for the font dataset when $s = 10$. Even if only fonts are detected (a), a large proportion of users are almost unique (i.e., have anonymity set size at most three). However, adding available screen size (b) greatly increases uniqueness and vanishes large anonymity sets that were visible before.

age set size drops fast; we observed 135.15 for $s = 10$ in case of the font dataset, 58.21 and 4344.93 for the app and location dataset, respectively. In case of the appli-



(a)



(b)

Fig. 10. Targeted fingerprinting: uniqueness in the font dataset with the size of anonymity sets when $s = 5$.

cation dataset, the average anonymity set size is 1.96 at $s = 29$. In case of the other datasets, the average never falls below 3, and we found 3.99 at $s = 76$ for the app dataset, and 3.96 at $s = 73$ for the font dataset.

While these results may not be promising to the adversary, measuring the proportion of unique users, whose anonymity sets are singletons, provides a different view. Unique users do not share the same fingerprint value with any other users in the dataset. As Fig. 12 shows, the proportion of unique users remains relatively low. The proportion does not change too much even for almost unique users which have anonymity set with size at most 3.

As regards linking attack against TBB users, we find that the available screen size as an auxiliary information (in addition to the detection of 10 fonts) makes this attack a significant threat, as 8,303 users (19.02%) are unique in this case. Without available screen size, less than 0.37% of users are unique with $s = 10$, but 7.89% of users are almost unique with $s = 25$ (Fig. 14; with available screen size in Fig. 16). We provide further details in Fig. 13 on the apps dataset, where uniqueness

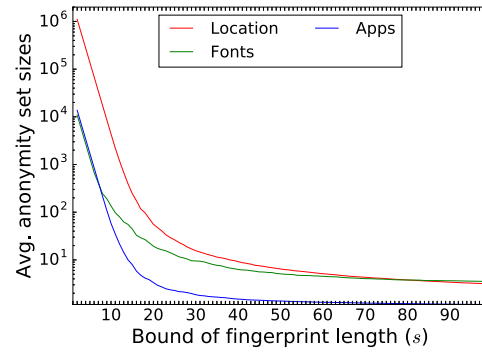


Fig. 11. General fingerprinting: average size of anonymity sets as the function of s .

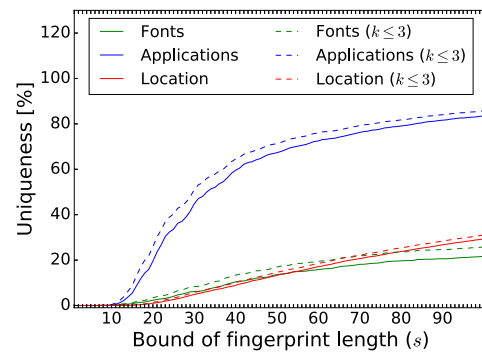


Fig. 12. General fingerprinting: proportion of unique and almost unique ($k \leq 3$) users as the function of s .

was the highest, that is, 76.30% of users have a size of anonymity set smaller than 4. This means that even if all 50 apps need to be pre-defined for all users, which is the case for iOS 9 at present, privacy threats still persist. For the location dataset, we find that 29% of all users are unique, and the average size of the anonymity set is 3.13 with $s = 100$ (Fig. 15). That is, an authority would need to deploy at least 100 cameras in the city to re-identify every third person on average in a pseudonymized CDR dataset (see Section 5.3).

7 Conclusion

In this paper we addressed two fingerprinting problems with constraints on the maximum size of the fingerprint. Targeted fingerprinting is executed on a specific user, whose profile and identity are known to the adversary, and the goal is to re-identify him/her in an anonymous context, where only a limited number of his

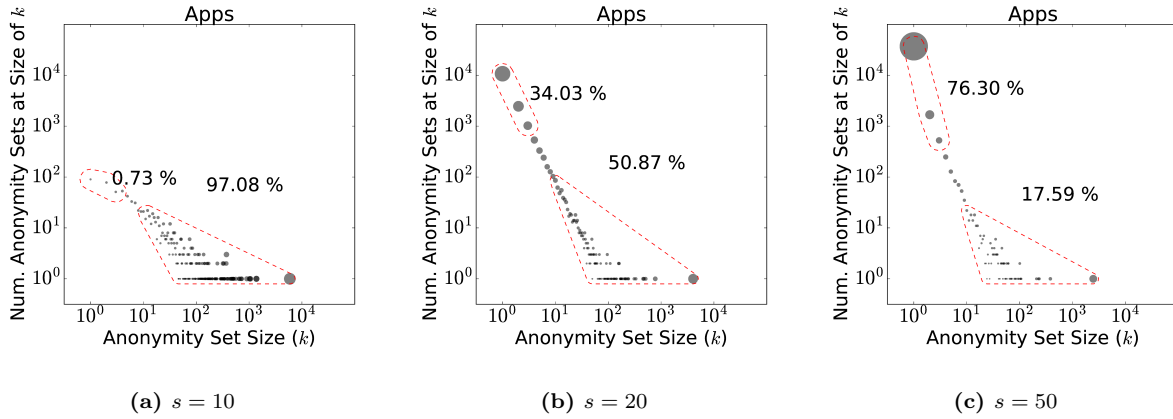


Fig. 13. General fingerprinting: distribution of size of anonymity sets for different values of s on the apps dataset.

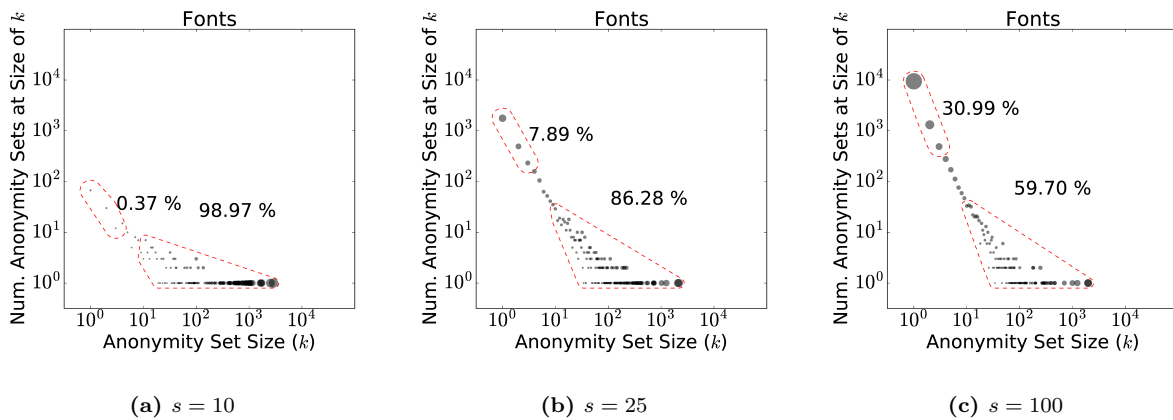


Fig. 14. General fingerprinting: distribution of size of anonymity sets for different values of s on the font dataset.

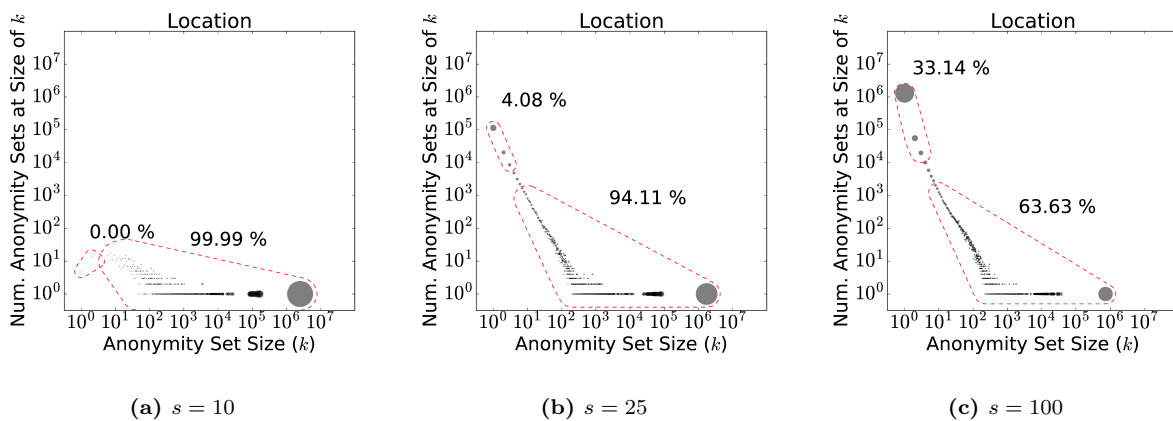


Fig. 15. General fingerprinting: distribution of size of anonymity sets for different values of s on the location dataset.

attributes can be queried. General fingerprinting is used to link the activities of multiple users without their re-identification in an anonymous context, where only a limited number of their attributes can be queried. We

provided essentially the best fingerprinting algorithms for both problems.

We evaluated these fingerprinting methods on real-life examples focusing on web tracking, smartphone user identification and mass de-anonymization of location

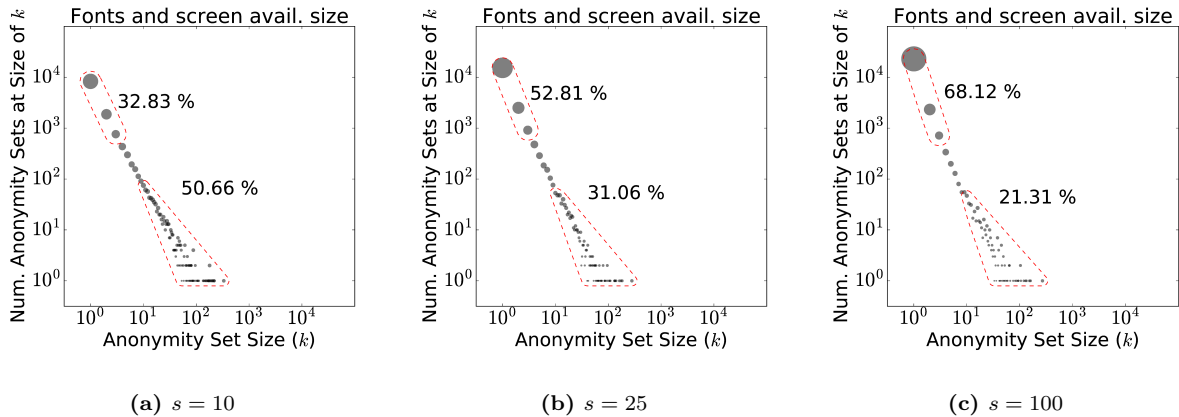


Fig. 16. General fingerprinting: distribution of size of anonymity sets for different values of s on the font dataset (with screen resolution and available screen size).

datasets. We showed that targeted fingerprinting is a significant threat to user privacy, as a significant ratio of users can be re-identified with our technique in case of web tracking and smartphone identification. This holds even when tighter constraints are applied than the default. Our results with general fingerprinting also support this claim; we found that 19.02% of web browsers could be uniquely identified by combining the presence of 10 fonts and the available screen size. We also found that by installing 100 surveillance cameras in a city of 4 million people, one can re-identify on average every fourth person in pseudonymized CDR dataset. Moreover, by querying 50 apps on a smartphone, which can be even pre-defined and identical for all users, one can obtain a unique fingerprint value for more than 65% of users in a dataset of 55k users.

All these findings show a conceptual weakness of limiting attribute access as a privacy guard.

8 Acknowledgement

The authors are thankful for Viktor Fábíán and Jagdish Prasad Achara for the discussions and their advises related to the iOS platform. This work was carried out during the tenure of an ERCIM 'Alain Bensoussan' Fellowship Programme. Some icons we used in the illustrations were made by Freepik (from www.flaticon.com)

References

[1] European Commission. Proposal for European Parliament

and the Council (General Data Protection Regulation), 2012. <http://eur-lex.europa.eu/LexUriServ/LexUriServ.do?uri=COM:2012:0011:FIN:EN:PDF>.

- [2] Apple uikit framework reference, entry on `canOpenURL` call. https://developer.apple.com/library/ios/documentation/UIKit/Reference/UIApplication_Class/#//apple_ref/occ/instm/UIApplication/canOpenURL., 2016.
- [3] G. Acar, C. Eubank, S. Englehardt, M. Juarez, A. Narayanan, and C. Diaz. The web never forgets: Persistent tracking mechanisms in the wild. In *ACM CCS*, pages 674–689, 2014.
- [4] J. P. Achara, G. Acs, and C. Castelluccia. On the unicity of smartphone applications. In *Proceedings of WPES*, pages 27–36. ACM, 2015.
- [5] Article 29 Data Protection Working Party. Opinion 05/2014 on anonymization techniques, April 2014.
- [6] K. Boda, A. M. Földes, G. G. Gulyás, and S. Imre. User tracking on the web via cross-browser fingerprinting. In P. Laud, editor, *Information Security Technology for Applications*, volume 7161 of *LNCS*, pages 31–46. 2012.
- [7] T. Bujlow, V. Carela-Espanol, J. Sole-Pareta, and P. Barlet-Ros. Web tracking: Mechanisms, implications, and defenses. In <http://arxiv.org/abs/1507.07872>, 2015.
- [8] X. Cai, X. C. Zhang, B. Joshi, and R. Johnson. Touching from a distance: Website fingerprinting attacks and defenses. In *ACM CCS*, 2012.
- [9] C. Chekuri and A. Kumar. Maximum coverage problem with group budget constraints and applications. In *Approximation, Randomization, and Combinatorial Optimization*, pages 72–83. Springer LNCS, 2004.
- [10] Y.-A. de Montjoye, C. A. Hidalgo, M. Verleysen, and V. D. Blondel. Unique in the crowd: The privacy bounds of human mobility. *Scientific Reports, Nature*, March 2013.
- [11] Y.-A. de Montjoye, L. Radaelli, V. K. Singh, and A. Pentland. Unique in the shopping mall: On the reidentifiability of credit card metadata. *Science*, 347(6221), January 2015.
- [12] P. Eckersley. How unique is your web browser? In *PETS*, pages 1–18, 2010.
- [13] U. Feige. A threshold of $\ln(n)$ for approximating set cover. *Journal of the ACM*, 45(4):634–652, July 1998.

- [14] gacar. Improve persistence and webfont compatibility of font patch, comment #13. <https://trac.torproject.org/projects/tor/ticket/5798comment:#13>, 2013.
- [15] M. Graham and S. D. Sabbata. The anonymous internet. <http://geography.oii.ox.ac.uk/?page=tor>, 2014.
- [16] G. G. Gulyás, G. Acs, and C. Castelluccia. Update your tor browser settings - otherwise it is less anonymous than you would think. <https://gulyas.info/blog/read/16/2015-11-30-Update-your-TOR-Browser-settings-otherwise-it-is-less-anonymous-than-you-would-think.php>, 11 2015.
- [17] J. Hayes and G. Danezis. k-fingerprinting: a robust scalable website fingerprinting technique. In <http://arxiv.org/abs/1509.00789>, 2015.
- [18] D. Lemire, L. Boytsov, and N. Kurz. SIMD compression and the intersection of sorted integers. *CoRR*, abs/1401.6399, 2014.
- [19] J. Marshall. Twitter is tracking users' installed apps for ad targeting. <http://blogs.wsj.com/cmo/2014/11/26/twitter-is-tracking-users-installed-apps-for-ad-targeting/>, 11 2014.
- [20] E. C. Mike Perry and S. Murdoch. The design and implementation of the tor browser [draft]. <https://www.torproject.org/projects/torbrowser/design/>, 5 2015.
- [21] R. Motwani and Y. Xu. Efficient algorithms for masking and finding quasi-identifiers. In *VLDB*, 2007.
- [22] G. Nemhauser, L. Wolsey, and M. Fisher. An analysis of approximations for maximizing submodular set functions i. *Mathematical Programming*, 14(1):265–294, 1978.
- [23] N. Nikiforakis, A. Kapravelos, W. Joosen, C. Kruegel, F. Piessens, and G. Vigna. Cookieless monster: Exploring the ecosystem of web-based device fingerprinting. In *IEEE Symposium on S&P*, pages 541–555, May 2013.
- [24] L. Olejnik, C. Castelluccia, and A. Janc. On the uniqueness of web browsing history patterns. *Annals of Telecommunications*, 69(1), February 2014.
- [25] A. J. Oliner, A. P. Iyer, I. Stoica, E. Lagerspetz, and S. Tarkoma. Carat: Collaborative energy diagnosis for mobile devices. In *ACM SenSys*, 2013.
- [26] H. T. T. Truong, E. Lagerspetz, P. Nurmi, A. J. Oliner, S. Tarkoma, N. Asokan, and S. Bhattacharya. The company you keep: Mobile malware infection rates and inexpensive risk indicators. In *WWW*, 2014.
- [27] H. Zang and J. Bolot. Anonymization of location data does not work: A large-scale measurement study. In *MobiCom*, 2011.