



HAL
open science

AP-Attack: A Novel User Re-identification Attack On Mobility Datasets.

Mohamed Maouche, Sonia Ben Mokhtar, Sara Bouchenak

► **To cite this version:**

Mohamed Maouche, Sonia Ben Mokhtar, Sara Bouchenak. AP-Attack: A Novel User Re-identification Attack On Mobility Datasets.. MobiQuitous 2017 - 14th EAI International Conference on Mobile and Ubiquitous Systems: Computing, Networking and Services, Nov 2017, Melbourne, Australia. pp.48-57, 10.1145/3144457.3144494 . hal-01785155

HAL Id: hal-01785155

<https://hal.science/hal-01785155>

Submitted on 22 Jun 2018

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

AP-Attack: A Novel User Re-identification Attack On Mobility Datasets

Mohamed Maouche
Universite de Lyon, CNRS
INSA Lyon, LIRIS, UMR5250
F69622, France
mohamed.maouche@insa-lyon.fr

Sonia Ben Mokhtar
Universite de Lyon, CNRS
INSA Lyon, LIRIS, UMR5250
F69622, France
sonia.benmokhtar@insa-lyon.fr

Sara Bouchenak
Universite de Lyon, CNRS
INSA Lyon, LIRIS, UMR5250
F69622, France
sara.bouchenak@insa-lyon.fr

ABSTRACT

Since the advent of hand held devices (e.g., smartphones, tablets, smart watches) with Ubiquitous computing and the wide popularity of location-based mobile applications, the amount of captured user location data is dramatically increasing. However, the gathering and exploitation of this data by mobile application providers raises many privacy threats as sensitive information can be inferred from it (e.g., home and work locations, religious beliefs, sexual orientations and social relationships). To address this issue a number of data obfuscation techniques (also called Location Privacy Protection Mechanisms or LPPMs) have been proposed in the literature. One of the existing methods to assess the effectiveness of LPPMs is to test them against user re-identification attacks. The aim of these attacks is to break user anonymity by re-associating data obfuscated using a given LPPM with user profiles built from user past mobility. In this paper, we present AP-Attack a novel re-identification attack that relies on a heatmap representation of user mobility data. Our experiments run against three representative LPPMs of the literature using four real mobility datasets show that AP-Attack succeeds in re-identifying up to 79% users in non-obfuscated data, +27% more users than POI-Attack and PIT-Attack two well known state-of-the-art attacks. We also present a simple technique to improve user protection against our attack, which relies on a user-centric application of multiple-LPPMs.

KEYWORDS

Security, Location Privacy, Mobility Trace, Re-identification attacks, Protection Mechanism

1 INTRODUCTION

With the raising number of mobile devices and the wide popularity of mobile applications, an increasing amount of mobility data is being gathered, processed and sometimes sold to third parties by application providers due to their inherent economic model. Examples of such applications include GPS navigation (e.g., Google Maps [21], Bing Maps [33]), location-based social networks (e.g., Swarm with Foursquare [13] or geo-gaming (e.g., Pokemon GO [37]). At the same time, following the open data movement, major socio-economic actors (e.g., telecommunication companies) and local authorities (e.g., cities) are pushed to give back their data to the society by publishing the datasets they are collecting about individuals [29] [46] [8]. However, as shown in various studies, the publication of user mobility data opens a number of privacy threats [27] [45] [47]. For instance, one can extract particular places where users regularly stop, also called Points Of Interest (POI) [16],

like the user's home location, work place [20], places of worship [14] or even discover the user health status if she regularly goes to the hospital. Moreover, by analyzing POIs of different users, social relationship can be discovered [5] and labels such as : siblings, colleagues, significant others... , can be associated to these relations.

To deal with this issue, various location privacy protection mechanisms (also called LPPMs) have been devised in the literature to protect the privacy of users when their mobility data is shared with applications. These mechanisms can be classified according to two usage scenarios. The first scenario, called the *online scenario* applies when users send their GPS coordinates to an application provider in order to get a geo-localized response (e.g., finding a restaurant in the user's vicinity, GPS navigation). In this context the LPPM, which runs on the client side can only act on the GPS coordinates sent by the user at a given time and place. Examples of such LPPMs include *Geo-Indistinguishability* [3], where Laplacian noise is added to each GPS coordinate, *CloakDroid* [32], where the GPS data is discretized using a grid or *Android Location Privacy Framework* [24], where various obfuscation techniques can be applied such as the generalization of a given location to the closest street, city, postal code and more. The second scenario, called the *offline scenario* applies when a given service provider collects a mobility dataset and needs obfuscation techniques to protect the participating users' privacy before releasing the dataset. In this context, the LPPM, which runs on the server side, has a broader view of the mobility of the overall population of users and can thus apply more sophisticated obfuscation techniques. Examples of such LPPMs include *GLOVE* [22], where mobility traces are merged together using a spatio-temporal similarity metric, *Never Walk Alone* [1] and its extension *W4M* [2], where cylindrical volumes wrap the movement of at least k different users together.

However, in both the online and the offline cases it is difficult to assess the effectiveness of the proposed LPPMs in practice. Indeed, LPPMs are generally evaluated either theoretically by proving the guarantees they offer to the users (e.g., k -anonymity [42] or differential privacy [11]) or practically by using custom privacy metrics that are often difficult to interpret, such as in [41] where POI retrieval metric is used to quantify privacy. Indeed, it is difficult to tell a data owner that aims at obfuscating her dataset whether obfuscating her dataset by enforcing k -anonymity with the *W4M* protocol [2] is better than obfuscating it by enforcing differential privacy with the *Geo-Indistinguishability* protocol [3]. A complementary way to assess the effectiveness of LPPMs is to rely on user re-identification attacks. Considering an obfuscated mobility dataset and a set of user profiles learnt from users past mobility, a

user re-identification attack tries to re-associate a portion of the obfuscated data to its originating user.

Literature contains a number of user re-identification attacks. These attacks can be distinguished according to two key elements: the user profiles they build from users past mobility and the distance metric they use to compare obfuscated data with user profiles. In this paper, we chose two state-of-the-art attacks POI-Attack [40] and PIT-Attack [15]. In the former, a user profile is represented by the list of POIs visited by the user while in the latter a user profile is represented by a Markov chain between the POIs visited by the user. However, none of these attacks consider the past mobility of users as a whole (i.e., considering both the places where the users stop and the trajectories that lead to these places).

In this paper, we first present AP-Attack (All Points Attack) a novel attack in which a user profile is represented as a heat-map. We compare the performance of AP-Attack to the two above attacks on four real mobility datasets and show that AP-Attack succeeds in re-identifying up to +27% more users than POI-Attack and up to +34% more users than PIT-Attack, reaching a re-identification rate of up to 79% in non-obfuscated data. We further use AP-Attack in addition to state-of-the-art attacks to show the lack of resilience of three state-of-the-art LPPMs (i.e., Geo-I [3], Promesse [41] and W4M [2]) to protect the users of four real datasets. Results show that none of the studied LPPMs succeeds in protecting all the users. Secondly, we study the vulnerability of individual users to re-identification attacks when their data is obfuscated using the above three LPPMs. Results show that users are not equal in front of re-identification attacks, some can not be protected by the considered LPPMs, some are naturally protected against the attacks, while others can be protected by one or multiple LPPMs. Using this observation, that to the best of our knowledge we are the first to establish, we propose a Multi-LPPM user-centric obfuscation technique, which outperforms all the evaluated LPPMs.

The remaining of this paper is structured as follows. First, we present in Section 2, a background on location privacy. Then, we present in Section 3, a model for re-identification attacks and a new re-identification attack AP-Attack. Further in Section 4, we evaluate AP-Attack and two state-of-the-art attacks of the literature POI-Attack and PIT-Attack against state-of-the-art LPPMs using four real datasets. As well as a technique to improve user protection using a Multi-LPPM approach based on a user-centric analysis. Finally we describe related research works in Section 5 before concluding the paper in Section 6.

2 BACKGROUND

We present in this section a set of background definitions related to mobility traces (Section 2.1) and to Location Privacy Protection Mechanisms (Section 2.2).

2.1 Mobility Traces

A mobility trace is constituted of a sequence of spatio-temporal points (lat, lng, t) associated to a given user, where lat and lng correspond to the latitude and longitude of GPS coordinates while t is a time stamp. The top part of Figure 1 shows a visual representation of a mobility trace (spacial elements only) of a given user collected in the city of San Francisco.

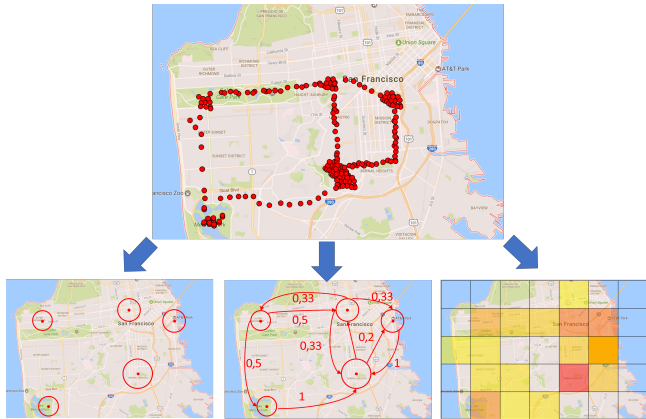
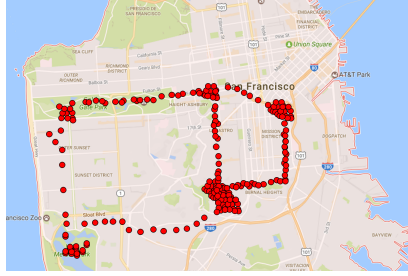


Figure 1: Various representations of mobility traces

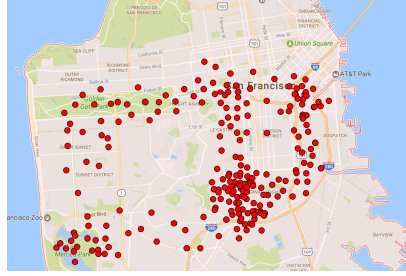
In order to associate semantic information to user raw mobility traces, various mobility models can be built from these traces. The bottom part of Figure 1 shows examples of such models. In this figure, the left part represents a mobility model in which only user points of interest (POIs) have been extracted. POIs are particular places where a user has stopped for a given amount of time. They are extracted from raw traces using spatio-temporal clustering algorithms such as [49] [23]. POIs may reveal personal information such as a user’s home place, work place or even sexual orientation and religious beliefs. The central part of the figure represents a mobility model in the form of a Markov chain between user POIs. This model is richer than the former one as it captures user mobility habits between POIs (e.g., the probability that the user goes to her favorite Japanese restaurant after going to the movie theatre). Finally, the right part of the figure represents a mobility model in which the map has been split into cells and the raw data has been projected into these cells in the form of a heat-map. Specifically, in this model, the intensity of the color of a given cell is relative to the frequency of user visits in the corresponding area of the map. Even though this model does not convey detailed temporal information about the user mobility, it is the only one to capture information about user trajectories. We will later use this model to build a novel user re-identification attack presented in Section 3.2.

2.2 Location Privacy Protection Mechanisms - LPPMs

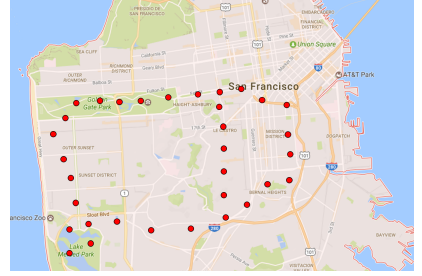
To overcome the threats affecting location privacy, Location Privacy Protection mechanisms (LPPMs) have been proposed in the literature. LPPMs generally take as input a mobility trace (sometimes composed of a single record [25]) or a set of mobility traces and alter these traces in order to produce obfuscated traces. LPPMs can be used in an online fashion, where each record is obfuscated before being sent to an application provider, or offline, where all the traces will be obfuscated at once. Furthermore, LPPMs are often classified depending on the privacy guarantees they offer to the users. There exist two major privacy guarantees presented in the literature: k -anonymity [42] and differential privacy [11]. The k -anonymity property states that a user is hidden among a set of $k - 1$ other users with similar properties. In the context of mobility data this



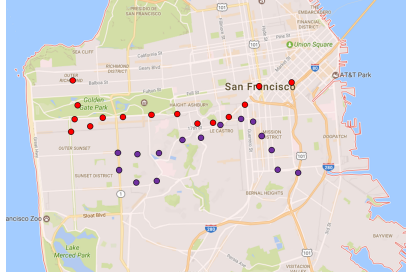
(a) One non-obfuscated trace



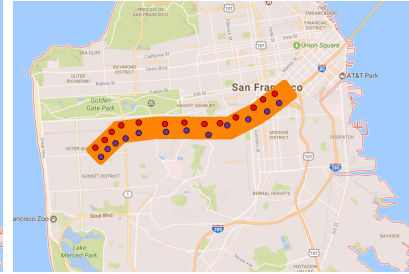
(b) Example of Geo-I applied to the trace in figure 2a



(c) Example of Promesse applied to the trace in figure 2a



(d) Two non-obfuscated traces



(e) Example of W4M applied to the traces in figure 2d

Figure 2: Illustration of LPPMs applied to mobility traces

translates to the ability to hide a given user in a geographical zone (called a cloaking area) where there are at least $k - 1$ other users [4]. Among the LPPMs that enforce k -anonymity, CliqueCloak [17] use a trusted third party to compute cloaking areas, PRIVE [19] has the same principle but relies on peer-to-peer communication between users to compute the cloaking areas. These two LPPMs allow the protection of a given geo-located point (i.e., online scenario) but do not consider a mobility trace as a whole. Instead, Wait 4 Me (W4M) [2] allows to enforce k -anonymity on mobility traces by extending k -anonymity to (k, δ) -anonymity. In this context, a user mobility trace will be hidden within $k - 1$ traces inside a cylindrical volume of radius δ . Figure 2e shows the application of W4M on the two mobility traces of Figure 2d. From these two figures, we observe that the two traces have been distorted to fit into the same cylindrical zone.

Differential privacy [11], which has initially been proposed for database systems, ensures that the result of an aggregate query over a table should not be significantly affected by the presence or absence of one single element of this table. This concept has been adapted to mobility data in an LPPM called Geo-Indistinguishability (Geo-I) [3]. In Geo-I, differential privacy is ensured by adding spatial noise to location data generated using a two dimensional Laplacian distribution. An example of applying Geo-I to a mobility trace of Figure 2a is depicted in Figure 2b. In this figure, we observe that each point in the original trace has been translated due to the added noise. As such, it is more difficult to infer information such as user POIs.

In addition to the above LPPMs, there exist other LPPMs that

try to protect user mobility traces by removing significant information from the traces such as user POIs. Among these LPPMs, Promesse [41] reaches this objective by distorting the temporal dimension of the mobility trace. Specifically, Promesse erases user POIs by using a speed smoothing technique, which assures that between each successive points in the obfuscated trace the distance and time difference are the same. An example of applying Promesse to a mobility trace of Figure 2a is depicted in Figure 2c. In this figure, we observe that POIs have been removed yet it is still possible to reason about user trajectories.

While the above LPPMs offer various theoretical or practical guarantees to protect the privacy of the users, it is difficult to guarantee resilience against powerful re-identification attack with background knowledge. In this paper, we show how re-identification attacks are able to break through the protection of state-of-the-art LPPM with different theoretical and practical guarantees.

3 AP-ATTACK: MODEL AND ALGORITHM

User re-identification attacks aim at linking user obfuscated data to her former mobility data. It is worth mentioning that the terminology *de-anonymization* can be found in place of *re-identification* (e.g., [15]). We would rather use *de-anonymization* to describe the process of finding a user real identity (e.g., name, address...) while re-identification describe the process of recovering a user ID in the system.

In the following, we present a general model for re-identification attacks Section 3.1, before describing our novel re-identification attack in Section 3.2.

3.1 Modelling Re-Identification Attacks

Let $U = \{U_1, U_2, \dots, U_N\}$ be the set of users in the system. The first phase of a re-identification attack is the training phase in which the adversary builds a knowledge base about the users in the system. In real systems, this phase may correspond to a period of time where users were using a geo-located service without protecting their mobility data. This phase is depicted in the left part of Figure 3. Specifically, we assume that for each user U_i , the adversary has access to a set of mobility traces corresponding to her past mobility, i.e., KD_i (where KD stands for Known user Data). Specifically, the set of all mobility traces known by the adversary is noted $KD = \{KD_1, KD_2, \dots, KD_n\}$. From each of these traces KD_i , we assume that the adversary builds a user profile $p(KD_i)$ that characterizes the user mobility as depicted in the bottom left part of Figure 3 (Step (1)). This profile is specific to each re-identification attack as further discussed in Section 4.1.

The second phase of a re-identification attack is depicted in the right part of Figure 3. In this phase, we assume that the adversary obtained a set of anonymous mobility traces (Step (2) in the figure), i.e., $UD = \{UD_1, UD_2, \dots, UD_m\}$ (where UD stands for Unknown user Data). Then, from each anonymous trace UD_i , the adversary builds a profile $p(UD_i)$ containing important information of the trace (Step (3) in the figure). Finally, a re-identification attack **A** run by the adversary tries to re-associate each extracted profile $p(UD_i)$ with profiles of known users, i.e.,

$$\begin{aligned} A : UD &\rightarrow U \\ UD_i &\mapsto A(UD_i, KD) = U_a \end{aligned}$$

A key element for the success of a re-identification attack is the similarity metric used to compare anonymous data with known user profiles (Step (4) in Figure 3). In addition to the way user profiles are modelled, the similarity metric is the second element, which is specific to each re-identification attack. If many anonymous traces are given as input to a re-identification attack, the attack is iterated on each element of UD_i as depicted in Algorithm 1. The success of an attack is then computed based on the number of correct re-associations the attack performs between anonymous traces and known user profiles (See line 9). To do this, we employ an oracle Id able to disclose for each anonymous trace UD_i its owner identity $Id(UD_i) = u(UD_i)$. This way, we can compute the **user re-identification rate** (Equation 1).

$$r(A_k, KD, UD) = \frac{\sum_{UD_i} \begin{cases} 1 & \text{If } A_k(UD_i, KD) = Id(UD_i) \\ 0 & \text{Else} \end{cases}}{|UD|} \quad (1)$$

3.2 AP-Attack Design Principles

We present in this section **AP-Attack** (All Points Attack) a novel re-identification attack that uses the whole user mobility data to form user profiles. Specifically, instead of focusing on a sub-set of points (e.g., those constituting POIs), AP-Attack aggregates all the points enclosed in a user mobility trace into a heat-map structure. More precisely, as shown in Figure 4, the map is subdivided into a

Algorithm 1 Re-identification attack

```

1: function  $A(UD, KD)$ 
2:    $UP \leftarrow \{p(UD_i) \setminus \forall i\}$ 
3:    $KP \leftarrow \{p(KD_j) \setminus \forall j\}$ 
4:    $matches \leftarrow \emptyset$ 
5:   for  $i \leftarrow [1, |UD|]$  do
6:      $j \leftarrow \arg \min_{0 < j \leq |KD|} d(UP_i, KP_j)$ 
7:      $matches \leftarrow matches \cup (Id(UD_i), Id(KD_j))$ 
8:   end for
9:    $rate \leftarrow \frac{\sum_{(u', u) \in matches} \begin{cases} 1 & \text{if } u' = u \\ 0 & \text{else} \end{cases}}{|UD|}$ 
10:  return  $(rate, matches)$ 
11: end function

```

grid with cells of the same size. Then, in each cell the number of records found in it is computed. As such, each cell will reflect the intensity of user movement in the corresponding geographical zone. This allows distinguishing between extremely, moderately, slightly frequented cells and unfrequented cells. Thereby, $p_{AP}(KD_i)$ return a probability distribution where each value $p_{AP}(KD_i)^{(k)}$ represents the probability that the owner of the trace U_i goes through the cell k . In order to be able to take into consideration the whole world map, the representation of each heatmap is a mapping between unique cells that the user passed by and their probability. This way each user would have a different sized map adapted to how wide her mobility was. This can also be seen as a sparse matrix.

Furthermore, we translate the distance between two profiles with the distance between two probability distributions. To compute this distance we can rely on classical distance metrics between probability distributions such as the ones surveyed in [7]. With respect to the experiments we did, one of the best metric to choose from is the Topsoe divergence defined in Equation 2. Where P and Q represent the list of cells in the two heat-maps we compare. So P_i is the probability of the user represented by the heat-map P going through the i th cell. This divergence is based on Shanon's concept of probabilistic uncertainty or entropy. It is a derived symmetric version of the Kullback Leibler divergence [7] which measures the information deviation. This is adapted to our case since we measure how much a heat-map can be used to characterize the mobility of a user that is represented by an other heat-map.

$$d_{Topsoe}(P, Q) = \sum_i \left[P_i \ln \left(\frac{2P_i}{P_i + Q_i} \right) + Q_i \ln \left(\frac{2Q_i}{P_i + Q_i} \right) \right] \quad (2)$$

In order to re-identify an anonymous trace UD_i , we match the trace with U_j one of the user of U whose trace KD_j minimize $d(p_{AP}(UD_i), p_{AP}(KD_j))$.

$$ie : A_{AP}(UD_i, KD) = \arg \min_{KD_j \in KD} (d(p_{AP}(UD_i), p_{AP}(KD_j)))$$

This new attack that not only take into consideration the POIs but also the mobility as a hole provides an effective counter-measure against LPPMs that are based on erasing POIs to raise the privacy level of a mobility trace.

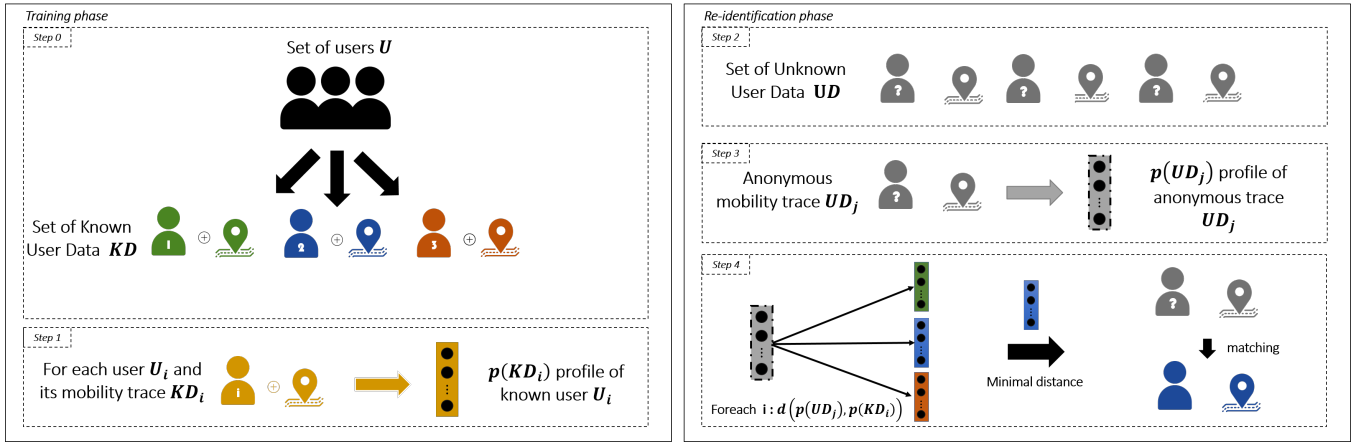


Figure 3: Re-identification attacks process from collecting phase to re-identification phase

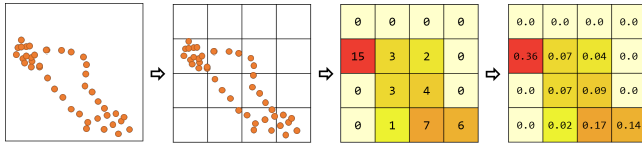


Figure 4: From mobility trace to user profile in AP-Attack

4 EVALUATION

We present in this section the evaluation of AP-Attack, we start by presenting the attacks and LPPMs used in this evaluation and how they have been configured (Section 4.1 to 4.3), our used datasets (Section 4.4) and our experimental setup (Section 4.5). Then, we present the performance of our proposed AP-Attack compared to state-of-the-art attacks (Section 4.6). We then demonstrate the lack of resilience of three representative LPPMs of the literature (Section 4.7). Finally, we use a user-centric approach in order to apply a Multi-LPPM obfuscation technique (Section 4.8).

4.1 Competitors

In this section we describe POI-Attack [40] and PIT-Attack [15] two state of the art attacks against which we compare the performance of AP-Attack.

4.1.1 Points Of Interest Attack - POI-Attack. This attack uses Points of interest (POIs) to characterize users' profiles. Therefore $p_{poi}(KD_i)$ is the set of POIs extracted from the trace KD_i . Those points are extracted using clustering algorithms such as the ones presented in [49] [23] parameterized with the diameter of a geographical zone where a user has stopped and a minimum duration characterizing her stop. To measure the similarity between two sets of POIs, each POI of the first set is associated with the geographically closest POIs in the second set. The dissimilarity between the two sets will be equal to the median of all the geographical distances, which is computed as follows in the Equation 3. Where X and Y are the sets of POIs for each trace and $d(X_r, Y_t)$ computes the geographical distance between two POIs X_r and Y_t .

$$d_{POISets}(X, Y) = \text{median} \left[\left\{ \min_t [d(X_r, Y_t)] \setminus \forall r \right\} \cup \left\{ \min_r [d(X_r, Y_t)] \setminus \forall t \right\} \right] \quad (3)$$

4.1.2 Probabilistic Inter-POI Transition Attack - PIT-Attack [15]. In addition to extracting POIs, this attack takes into consideration the transition probability from one POI to another. Specifically, the authors rely on mobility Markov chains [16] where the states are POIs ($P = P_1, P_2, \dots, P_k$) ordered by the number points in each POI and the edges' labels are transitions probabilities between POIs (t_{P_i, P_j}). This is done by computing the proportion of transition between each POI in the mobility traces. In order to compute the distance between two mobility Markov chains, two informations are taken into account : the geographical distance between POIs and the weight of each POI. The weight of a POI is computed using the proportion of points contained inside the POI. More precisely the authors proposed many distance metrics to compare Markov chains. The most effective one is the *stats-prox* distance which is a combination of two distances: the stationary distance and the proximity distance (Equation 4). The stationary distance (Equation 5) sums the weighted geographical distances between each combination of two POIs if the distance is lower then a parameter d_0 . And the proximity distance (Equation 6) after ranking the POIs by their weight in each Markov Chain. It adds scores r_i if two POIs of the rank i are closer than a parameter Δ . The score is halved after each rank $r_i = \frac{1}{2}r_{i-1}$ and r_0 is a parameter. The dissimilarity between the two Markov chain is the inverse of the total score.

$$d_{stats-prox} \equiv \text{if}(d_{stat} > \gamma) d_{stat} \text{ else } d_{prox} \quad (4)$$

$$d_{stats}(P, Q) = \sum_{P_i, Q_j \in P \times Q} w(P_i) \times \begin{cases} d(P_i, Q_j) & \text{If } d(P_i, Q_j) < d_0 \\ 0 & \text{Else} \end{cases} \quad (5)$$

$$d_{prox}(P, Q) = \left(\sum_{i=1}^{\min(|P|, |Q|)} \begin{cases} r_i & \text{If } d(P_i, Q_i) < \Delta \\ 0 & \text{Else} \end{cases} \right)^{-1} \quad (6)$$

$$r_i = \frac{1}{2} r_{i-1}$$

This attacks rely almost exclusively on POIs, eliminating the information contained inside the trajectories. Also LPPMs focusing on the elimination of POIs yield to a inept attack as illustrated in Section 4.7.

4.2 Attacks Configuration

The three attacks evaluated in this paper AP-Attack, POI-Attack and PIT-Attack have a number of configuration parameters. Specifically, AP-Attack has a cell size parameter that we have fixed at 800 meters in this evaluation. After a number of calibration experiments, we have chosen this value because it was big enough to include POIs and was resilient to noisy traces. In addition, re-identification rates result for cells between 50 meters and 800 meters are approximately similar. Furthermore, POI-Attack and PIT-Attack require parameters for the extraction of the POIs from the traces. These parameters are the diameter of the clustering area (that we fixed at 200 meters) and the minimum time spent inside a POI (that we fixed at 1 hour). These values have been chosen after a series of experiments yielding to the best results. It is worth mentioning that in [40] POI-Attack was used in a different configuration. Indeed, the authors re-identified the obfuscated mobility traces against the non-obfuscated version of those traces, rather than using a different past mobility as a training knowledge. In consequence, the re-identification is easier.

4.3 LPPMs

To evaluate AP-Attack, we have chosen three representative LPPMs of the literature (see section 2.2): (1) Geo-I, which adds Laplacian noise to mobility traces and enforces a guarantee inspired from Differential privacy; (2) Promesse, which uses speed smoothing to erase POIs and (3) W4M, which alters traces to group them in cylindrical volumes hence enforcing k-anonymity. Each LPPM has a number of configuration parameters. These parameters have an impact on the privacy level offered to the users but also on the quality of the resulting obfuscated data. Due to a lack of space, we decided to configure each LPPM following a medium level of protection. This choice is motivated by the fact that our objective is not to find the best LPPM configuration but rather to show how with a reasonable alteration of the data. The LPPM do not succeed completely in protecting the user from re-identification. Other experiments with other configurations of the used LPPMs or using other LPPMs of the literature can be done using our available toolkit [31]. Specifically, Geo-I is configured with a parameter ϵ that has an impact on the amount of noise added to the data (the lower epsilon the higher the noise). We have fixed the value of this parameter to 0.01, which corresponds to a medium privacy level. Promesse is configured with a parameter α that corresponds to the distance between two successive sampling points. We have fixed this parameter to 200 meters. Finally, W4M is configured with two

Table 1: Description of datasets

Name	CabSpotting	Geolife	MDC	PrivaMov
# users	536	42	144	48
Localization	San Francisco	Beijing	Geneva	Lyon
# records	11.219.955	1.574.338	904.422	973.684

parameters, k representing the minimum number of users inside the cylindrical volume and the radius δ of the latter. We have fixed these parameters at $k = 2$ and $\delta = 600$ meters because W4M erases a lot of points making the dataset almost empty and those parameters guarantee privacy and availability of the data.

4.4 Datasets

We used four real mobility datasets in our experiments. These datasets are:(1) Cabspotting [39] that contains the mobility of 536 cab drivers in the city of San Francisco; (2) Geolife [48] that contains the mobility of 42 users mainly in the city of Beijing; (3) MDC [29] that contains the mobility data of 144 users in the city of Geneva and (4) PrivaMov [6] that contains the mobility of 48 students and staff members in the city of Lyon. To make the comparison fair between the datasets, we selected in each dataset the 30 most active successive days. We present in the table 1 a description of the datasets used in our experiments. The users are not active in all the days of the period some are more active than others. We consider as a mobility trace, the mobility of the user during all the period. In all the experiments described in this paper, we split the datasets into a period of 15 days used for the training phase and 15 days used for the re-identification phase. We run other experiments where the training and re-identification phases varied from 1 day to 23 days each to evaluate the impact of dataset splitting on re-identification attacks. We do not present these results in the paper due the lack of space, but the results are available in the companion technical report [31].

4.5 Experimental Setup

All of our experiment were carried out in a computer running an Ubuntu 14.04 OS with 50GB of RAM and 16 cores of 1.2Ghz each. Our testing application [31] written in Java & Scala and runs in the Java Virtual Machine 1.8.0.

4.6 Performance Of Re-identification Attacks

The first experiment we did was intended to compare the three considered re-identification attacks by measuring their re-identification rate on non-obfuscated data of the four considered datasets. Results are depicted in Figure 5. From this figure, we observe that AP-Attack outperforms the two other attacks on all the considered datasets. This experiment shows that sending mobility data "anonymously" (e.g., by using anonymous communication protocols such as TOR [10]) to application providers is not sufficient to protect the privacy of users as an adversary using re-identification attacks is able to recognize from 45% to 79% of the users in the four datasets. It is thus necessary for end users to rely on LPPMs to protect their data. From this experiment we also notice that Cabspotting is the

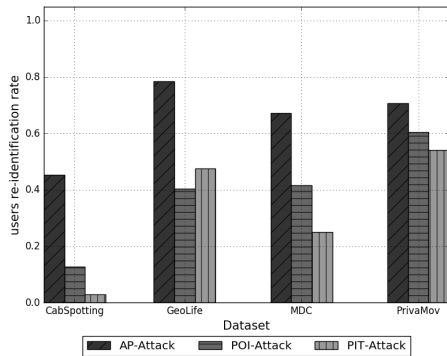


Figure 5: Performance of re-identification attacks

dataset where the users are the most intrinsically protected. This comes from the fact that Cab drivers have similar mobility patterns (e.g., they regularly go to the airport, famous hotels, malls and the taxi parking places). Instead, MDC, GeoLife and PrivaMov are related to users having different mobility habits, which makes them easier to re-identify.

4.7 LPPMs Effectiveness Against Re-identification Attacks

In this experiment, we compare the performance of the three considered LPPMs, i.e., Geo-I, Promesse and W4M. Specifically, we evaluate the re-identification rate obtained by the three former attacks on data obfuscated using these three LPPMs. Figure 6 shows the result of this experiment. In addition to the three LPPMs, we report the results obtained for non-obfuscated data, which we use as a baseline. At first glance, we observe the high level of privacy enforced by W4M in the PrivaMov dataset (11%) and by Promesse in the Cabspotting dataset (6%) against AP-Attack, which is the most successful attack. Nevertheless, these two LPPMs seem not to be sufficient to protect users in the GeoLife and MDC datasets where the re-identification rate reaches 48% and 36% for W4M and 68% and 46% for Promesse. We notice that the LPPMs that erase POIs as Promesse and W4M nullify the attack POI-Attack and PIT-Attack. For instance, in the Geolife Dataset for AP-Attack goes from 79% in the non-obfuscated data to 68% in the dataset protected with Promesse while PIT-Attack goes from 47% to 0%. Finally, we observe that Geo-I is the least efficient LPPM against re-identification attacks in the four datasets. We also notice that Geo-I affects less AP-Attack compared to POI-Attack. Indeed, AP-Attack goes down in average with -3% while POI-Attack goes down by -15% . The noise added to the points by Geo-I rarely gets them out of a cell, while the clustering algorithms used to form POIs suffer more from the noise. Summarizing, this experiment allows us to draw the following conclusions: (1) there is no one-size-fits-all LPPM, as the resilience of an LPPM to re-identification attacks depends on the underlying data; (2) users of a given dataset are not all equal in front of re-identification attacks, as on the four datasets there exist users that are never re-identified even in the absence of protection mechanisms (e.g., 54% for the best case with Cabspotting and 21% for the worst case with Geolife). These two observations motivate the

need of investigating multi-LPPM and user centric data obfuscation techniques as presented in the following section.

4.8 Improving Dataset Obfuscation

With regards to the results in the previous section, a user centric analysis is needed. We start this experiment by evaluating the sensitivity of individual users to re-identification attacks and show how a user can be mistaken for another (Section 4.8.1). We then investigate a multi-LPPM protection scheme (Section 4.8.2).

4.8.1 Sensitivity of users to re-identification attacks. This experiment shows the proportion of users protected by none, one or multiple LPPMs on each of our four datasets. In this experiment we used all the re-identification attacks (ie., AP-Attack, POI-Attack and PIT-Attack). Results are depicted in Figure 7. Overall these results allow us to draw the following conclusions: (1) there is a proportion of users that are not vulnerable to re-identification attacks (this proportion varies from 19% to 54% in the four datasets); (2) there is a proportion of users that can not be protected by the existing LPPMs in their current configuration (this proportion varies from 2% to 33% in the four datasets); (3) there is a proportion of users that can be protected by only one LPPM among those that we tested (this proportion varies from 19% to 42% in the four datasets) and (4) there is a proportion of users that can be protected by multiple LPPMs (this proportion varies from 12% to 37% in the four datasets). From these conclusions, which to the best of our knowledge we are the first to draw, it becomes natural to think of multi-LPPM obfuscation techniques as further discussed in the following section.

A user is mistakenly re-identified when her mobility in the anonymous trace is similar to the mobility of another user in the known dataset. We looked into the individual heat-maps in the datasets. We notice that often, the user share parts of his mobility in the anonymous trace with her known profile but also with another user and some little behavior changes make her mistakenly re-identified. While, we rarely found users that drastically change their behavior from known and unknown dataset. From this observation, we recommend for the design of LPPMs resilient against re-identification attacks, approaching the trace to wrong users in order to confuse the attacker rather than drastically altering the trace and degrading its utility for the service provider.

4.8.2 Towards Multi-LPPM Obfuscation. In this experiment, we decided to leverage the results obtained in the previous experiment to design a multi-LPPM obfuscation technique, the method is described in the Algorithm 2. Specifically, on each of our four datasets we built an obfuscated dataset as follows. For each user, we chose one mobility trace from the non-obfuscated dataset or one of the datasets obfuscated by an LPPM. To do this, we extracted the list of all LPPMs that successfully protected the user (line 4). Then, chose the trace according to a preference order (line 6). In consequence, the users that were insensitive to all the re-identification attacks. As they are naturally protected, it is better not to alter their corresponding portion of the data. Then, for the users that were protected by only one LPPM, we used the latter in our obfuscated dataset. Finally, for those users that were protected by more than one LPPM, we used in the following order Geo-I, Promesse or W4M

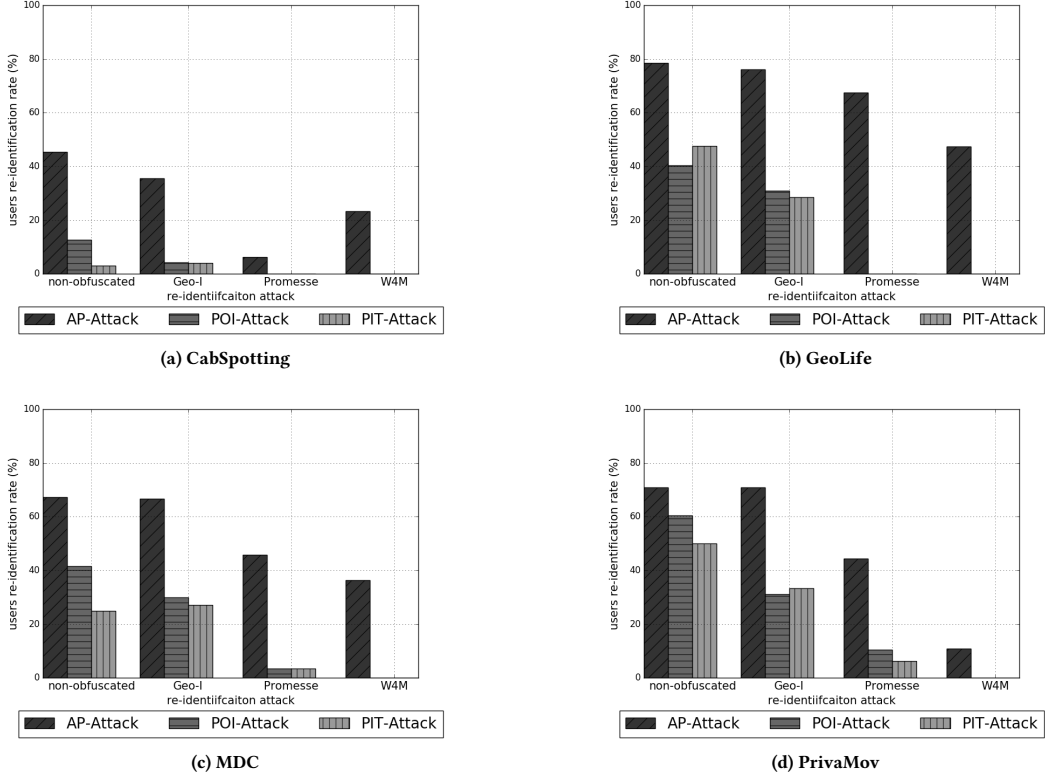


Figure 6: Performance of LPPMs

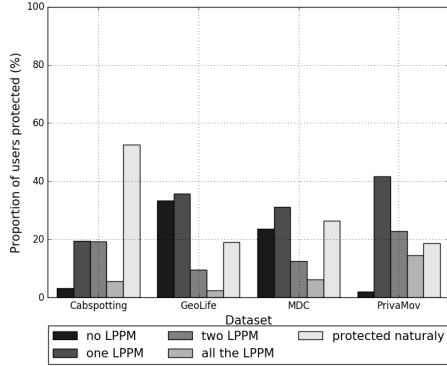


Figure 7: Proportion of users protected by a number of LPPM

to obfuscate their data. This choice was motivated by the degree of degradation obtained in the traces after obfuscation, which is lower when using Geo-I and Promesse than when using W4M.

The results depicted in Figure 8 show that our multi-LPPM and user-centric obfuscation technique called Hybrid in the figure outperforms all the existing LPPMs. Nevertheless, there are still users that are not protected by our multi-LPPM approach. This suggests that there is still room for proposing novel LPPMs. Our findings

Algorithm 2 Multi-LPPM user-centric dataset obfuscation

```

1: function HybridLPPM( $U, UD, ResultsAttacks$ )
2:    $dataset_{out} \leftarrow \emptyset$ 
3:   for all  $i \in U$  do
4:      $\ell \leftarrow getProtectingLPPMs(ResultsAttacks_i)$ 
5:     if  $|\ell| \neq 0$  then
6:        $sort(\ell)$ 
7:       // nonObfuscated  $\rightarrow$  GeoI  $\rightarrow$  Promesse  $\rightarrow$  W4M
8:        $t \leftarrow UD_i^{\ell[0]}$ 
9:       //  $UD^x$  is the dataset obfuscated with the LPPM  $x$ 
10:    else
11:       $t \leftarrow UD_i$ 
12:    end if
13:     $dataset_{out} \leftarrow dataset_{out} \cup t$ 
14:  end for
15: return  $dataset_{out}$ 
16: end function

```

suggest that efforts should go in the direction of a data-centric/user-centric approach. We showed in this section how a system designer or data owner going towards this direction by adapting the LPPMs and possibly their degree of protection according to the sensitivity of each user to re-identification attacks.

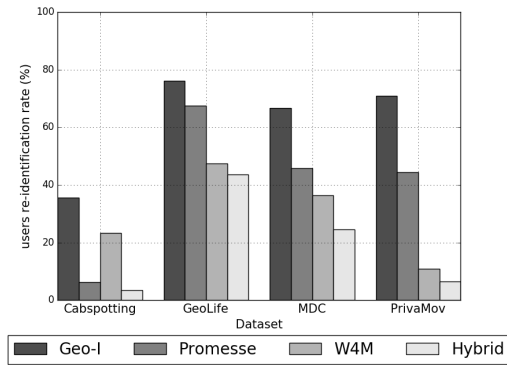


Figure 8: Performance of the Hybrid LPPM

5 RELATED WORK

The re-identification threat is affecting a wide variety of systems due to the wide scale gathering of user personal data by application providers. To measure this threat, a variety of re-identification attacks are being proposed in various systems such as Web-search systems [18] [38], face recognition [12], social networks [35] and recommender systems [36]. User re-identification attacks in the context of mobility data such as the one proposed in this paper, share a similar objective as the attacks above. Specifically, re-identification attacks demonstrate that the accumulation of mobility data about users allows the extraction of user profiles despite users hiding or changing their user ID. In order to protect against these threats various Location Privacy Protection Mechanism (LPPM) have been introduced. As previously discussed in this paper, LPPMs alter mobility traces in order to protect users against the inference of sensitive information about them. To evaluate the degree of protection offered by LPPMs to users, various privacy evaluation metrics have been used. Examples of such metrics include the POI retrieval rate [41, 43], which reflects the ability of a LPPM to hide user POIs in the obfuscated traces.

We find in the literature an increasing interest in location privacy. Indeed, as the work of De Montjoye & al. [9] have shown, the mobility of users acts as a fingerprint. They managed with only four points from the mobility traces to isolate a unique user with a 95% success rate. While uniqueness in the crowd is different than re-identification because it does not use new traces to compare against formerly gathered one, it still make proof of how a mobility trace discriminate users. Krumm et al. [28] managed to put a learning system which is able to label geographic places (Home, work, worship, shop...) with 73% success rate demonstrating the vulnerability of location privacy in semantic extraction from POIs. On the same line, Krumm [26] looked at user de-anonymization by finding users' home addresses, they were able to find users' homes by a median error of 60 meters but the white pages system they used was not effective enough to find user real identity. Ma & al. [30] studied another type of re-identification where the anonymous traces are intercalated between the records of the identified traces. Naini & al. [34] also used map grid to compare between users. They worked on a closed system and tried using a bipartite graph matching to associate users. Srivatsa & al [44] used social network as side

channel to re-identify users, they used a contact graph identifying meetings between users extracted from a set of traces and then used a correlation with a social network graph to match users mobility with their social network account.

6 CONCLUSION

In this paper, we presented a novel re-identification attack based on a heat-map representation of user profiles. We showed that this attack, which aggregates user mobility into a probability distribution acting as a fingerprint of user mobility, outperforms existing attacks on four real mobility datasets. Moreover, we studied the ability of three state-of-the-art LPPMs to protect users against re-identification attacks. The results showed that there is no one-size-fits-all LPPM. Instead, the degree of protection offered by LPPMs heavily depend on the underlying data. We then decided to further analyze how individual users are sensitive to re-identification attacks while being protected by various LPPMs. Our results have shown that users are not equal in front of re-identification attacks, some can not be protected by the considered LPPMs, some are naturally protected, while others can be protected by one or multiple LPPMs. This observation, that to the best of our knowledge we are the first to establish has lead us to the design of a Multi-LPPM user-centric obfuscation technique, which better resists re-identification attacks than existing LPPMs. Still, according to the considered datasets, we showed that a proportion of users can not be protected using this technique, which opens the door for future investigations in the field. In this context, we have shown that current state-of-the-art LPPM lack the resilience to protect user against LPPM in front of powerful re-identification attacks and demonstrated that data owners need to consider user-centric obfuscation for the protection of their datasets. Also, the LPPM affect the utility of the traces in different intensities. That's why, we need to investigate the Utility aspects of LPPMs and look for a good trade-off between Privacy and Utility.

REFERENCES

- [1] Osman Abul, Francesco Bonchi, and Mirco Nanni. 2008. Never Walk Alone: Uncertainty for Anonymity in Moving Objects Databases. In *Proceedings of the 2008 IEEE 24th International Conference on Data Engineering (ICDE '08)*. IEEE Computer Society, Washington, DC, USA, 376–385. <https://doi.org/10.1109/ICDE.2008.4497446>
- [2] Osman Abul, Francesco Bonchi, and Mirco Nanni. 2010. Anonymization of moving objects databases by clustering and perturbation. *Information Systems* 35, 8 (2010), 884–910. <https://doi.org/10.1016/j.is.2010.05.003>
- [3] Miguel E. Andrés, Nicolás E. Bordenabe, Konstantinos Chatzikokolakis, and Catuscia Palamidessi. 2013. Geo-Indistinguishability: Differential Privacy for Location-Based Systems. *Ccs'13 abs/1212.1* (2013), -. <https://doi.org/10.1145/2508859.2516735> arXiv:arXiv:1212.1984v2
- [4] Claudio Bettini, X Sean Wang, and Sushil Jajodia. 2005. Protecting Privacy Against Location-based Personal Identification. In *Proceedings of the Second VDLB International Conference on Secure Data Management (SDM '05)*. Springer-Verlag, Berlin, Heidelberg, 185–199. https://doi.org/10.1007/11552338_13
- [5] Igor Bilogrevic, Kevin Huguenin, Murtuza Jadhwal, Florent Lopez, Jean-Pierre Hubaux, Philip Ginzboorg, and Valteri Niemi. 2013. Inferring Social Ties in Academic Networks Using Short-Range Wireless Communications. *Wpes* (2013), 179–188. <https://doi.org/10.1145/2517840.2517842>
- [6] Antoine Boutet, Sonia Ben Mokhtar, and Vincent Primault. 2016. Uniqueness Assessment of Human Mobility on Multi-Sensor Datasets. (2016).
- [7] Sung-hyuk Cha. 2007. Comprehensive Survey on Distance / Similarity Measures between Probability Density Functions. *International Journal of Mathematical Models and Methods in Applied Sciences* 1, 4 (2007), 300–307. <https://doi.org/10.1007/s00167-009-0884-z>
- [8] Data team (UK's Cabinet Office). 2017. Open data in UK. (2017). <https://data.gov.uk/>

- [9] Yves-Alexandre de Montjoye, César A. Hidalgo, Michel Verleysen, and Vincent D. Blondel. 2013. Unique in the Crowd: The privacy bounds of human mobility. *Scientific reports* 3 (2013), 1376. <https://doi.org/10.1038/srep01376>
- [10] R Dingleline, N Mathewson, and P Syverson. 2004. Tor: The second-generation onion router. (2004).
- [11] Cynthia Dwork. 2008. *Differential Privacy: A Survey of Results*. Springer Berlin Heidelberg, Berlin, Heidelberg, 1–19. https://doi.org/10.1007/978-3-540-79228-4_1
- [12] M Farenzena, L Bazzani, A Perina, V Murino, and M Cristani. 2010. Person re-identification by symmetry-driven accumulation of local features. (2010), 2360–2367 pages. <https://doi.org/10.1109/CVPR.2010.5539926>
- [13] Foursquare-Labs. 2017. Swarm. (2017). <https://www.swarmapp.com>
- [14] L. Franceschi-Bicchierai. 2015. Redditor cracks anonymous data trove to pinpoint muslim cab drivers. (2015). <http://mashable.com/>
- [15] Sebastien Gams, Marc-Olivier Killijian, and Miguel Nunez del Prado Cortez. 2013. De-anonymization Attack on Geolocated Data. *2013 12th IEEE International Conference on Trust, Security and Privacy in Computing and Communications* (2013), 789–797. <https://doi.org/10.1109/TrustCom.2013.96>
- [16] Sébastien Gams, Marc-Olivier Killijian, and Miguel Nez Del Prado Cortez. 2011. Show Me How You Move and I Will Tell You Who You Are. *Transactions on Data Privacy* 4 (2011), 103–126. <https://doi.org/10.1145/1868470.1868479>
- [17] Bugra Gedik and Ling Liu. 2005. Location Privacy in Mobile Systems: A Personalized Anonymization Model. In *Proceedings of the 25th IEEE International Conference on Distributed Computing Systems (ICDCS '05)*. IEEE Computer Society, Washington, DC, USA, 620–629. <https://doi.org/10.1109/ICDCS.2005.48>
- [18] Arthur Gervais, Reza Shokri, Adish Singla, Srdjan Capkun, and Vincent Lenders. 2014. Quantifying Web-Search Privacy. *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security* (2014), 966–977. <https://doi.org/10.1145/2660267.2660367>
- [19] Gabriel Ghinita, Panos Kalnis, and Spiros Skiadopoulos. 2007. PRIVE: Anonymous Location-based Queries in Distributed Mobile Systems. In *Proceedings of the 16th International Conference on World Wide Web (WWW '07)*. ACM, New York, NY, USA, 371–380. <https://doi.org/10.1145/1242572.1242623>
- [20] Philippe Golle and Kurt Partridge. 2009. On the anonymity of home/work location pairs. In *International Conference on Pervasive Computing*. Springer, 390–397.
- [21] Google. 2017. Google maps. (2017). <https://maps.google.com>
- [22] Marco Gramaglia and Marco Fiore. 2015. Hiding Mobile Traffic Fingerprints with GLOVE. In *Proceedings of the 11th ACM Conference on Emerging Networking Experiments and Technologies (CoNEXT '15)*. ACM, New York, NY, USA, 26:1–26:13. <https://doi.org/10.1145/2716281.2836111>
- [23] Ramaswamy Hariharan and Kentaro Toyama. 2004. Project Lachesis: Parsing and Modeling Location Histories. In *Geographic Information Science: Third International Conference, GIScience 2004, Adelphi, MD, USA, October 20-23, 2004. Proceedings*, Max J Egenhofer, Christian Freksa, and Harvey J Miller (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 106–124. https://doi.org/10.1007/978-3-540-30231-5_8
- [24] B Henne, C Kater, M Smith, and M Brenner. 2013. Selective cloaking: Need-to-know for location-based apps. (2013), 19–26 pages. <https://doi.org/10.1109/PST.2013.6596032>
- [25] Christian S Jensen, Hua Lu, and Man Lung Yiu. 2009. Location Privacy Techniques In Client Server Architectures. *Privacy in Location-Based Applications* 5599 (2009), 31–58.
- [26] John Krumm. 2007. Inference Attacks on Location Tracks. *Pervasive Computing* 10, Pervasive (2007), 127–143. https://doi.org/10.1007/978-3-540-72037-9_8
- [27] John Krumm. 2009. A survey of computational location privacy. *Personal and Ubiquitous Computing* 13, 6 (2009), 391–399. <https://doi.org/10.1007/s00779-008-0212-5>
- [28] John Krumm and Dany Rouhana. 2013. Placer: Semantic Place Labels from Diary Data. In *Proceedings of the 2013 ACM International Joint Conference on Pervasive and Ubiquitous Computing (UbiComp '13)*. ACM, New York, NY, USA, 163–172. <https://doi.org/10.1145/2493432.2493504>
- [29] J K Laurila, Daniel Gatica-Perez, I Aad, Blom J., Olivier Bornet, Trinh-Minh-Tri Do, O Dousse, J Eberle, and M Miettinen. 2012. The Mobile Data Challenge: Big Data for Mobile Computing Research. In *Pervasive Computing*.
- [30] Chris Y T Ma, David K Y Yau, Nung Kwan Yip, and Nageswara S V Rao. 2013. Privacy vulnerability of published anonymous mobility traces. *IEEE/ACM Transactions on Networking* 21, 3 (2013), 720–733. <https://doi.org/10.1109/TNET.2012.2208983>
- [31] Mohamed Maouche. 2017. SFERA. (2017). <https://github.com/mmaouche-insa/SFERA/>
- [32] Kristopher Micinski, Philip Phelps, and Jeffrey S Foster. 2013. An Empirical Study of Location Truncation on Android. *Most '13* (2013). <http://www.mostconf.org/2013/papers/19.pdf>
- [33] Microsoft. 2017. Bing maps. (2017). <https://www.bing.com/maps>
- [34] F.M. Naini, J. Unnikrishnan, P. Thiran, and M. Vetterli. 2016. Where You Are Is Who You Are: User Identification by Matching Statistics. *IEEE Transactions on Information Forensics and Security* 11, 2 (2016), 358–372. <https://doi.org/10.1109/TIFS.2015.2498131> arXiv:1512.02896
- [35] A Narayanan, E Shi, and B I P Rubinstein. 2011. Link prediction by de-anonymization: How We Won the Kaggle Social Network Challenge. (2011), 1825–1834 pages. <https://doi.org/10.1109/IJCNN.2011.6033446>
- [36] Arvind Narayanan and Vitaly Shmatikov. 2008. Robust de-anonymization of large sparse datasets. *Proceedings - IEEE Symposium on Security and Privacy* (2008), 111–125. <https://doi.org/10.1109/SP.2008.33> arXiv:arXiv:cs/0610105v2
- [37] Niantic. 2017. Pokemon Go. (2017). <http://www.pokemongo.com>
- [38] Albin Petit, Thomas Cerqueus, Antoine Boutet, Sonia Ben Mokhtar, David Coquil, Lionel Brunie, and Harald Kosch. 2016. SimAttack: private web search under fire. *Journal of Internet Services and Applications* 7, 1 (2016), 2. <https://doi.org/10.1186/s13174-016-0044-x>
- [39] Michal Piorowski, Natasa Sarafijanovic-djukic, and Matthias Grossglauser. 2009. CRAW- DAD data set epfl/mobility (v. 2009-02-24). (2009). <http://crawdad.cs.dartmouth.edu/epfl/mobility>
- [40] Vincent Primault, Sonia Ben Mokhtar, Cédric Lauradoux, and Lionel Brunie. 2014. Differentially Private Location Privacy in Practice. *Most '14* October (2014). arXiv:1410.7744
- [41] Vincent Primault, Sonia Ben Mokhtar, Cédric Lauradoux, and Lionel Brunie. 2015. Time distortion anonymization for the publication of mobility data with high utility. *Proceedings - 14th IEEE International Conference on Trust, Security and Privacy in Computing and Communications, TrustCom 2015* 1 (2015), 539–546. <https://doi.org/10.1109/Trustcom.2015.417> arXiv:1507.0443
- [42] Pierangela Samarati and Latanya Sweeney. 1998. Generalizing Data to Provide Anonymity when Disclosing Information. In *Proceedings of the Seventeenth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS '98)*. ACM, New York, NY, USA, 188–. <https://doi.org/10.1145/275487.275508>
- [43] Reza Shokri, George Theodorakopoulos, Jean Yves Le Boudec, and Jean Pierre Hubaux. 2011. Quantifying location privacy. *Proceedings - IEEE Symposium on Security and Privacy* (2011), 247–262. <https://doi.org/10.1109/SP.2011.18>
- [44] Mudhakar Srivatsa and Mike Hicks. 2012. De-anonymizing Mobility Traces : Using Social Networks as a Side-Channel. *Proceedings of the 2012 ACM conference on Computer and communications security (CCS)* (2012), 628–637. <https://doi.org/10.1145/2382196.2382262>
- [45] Manolis Terrovitis. 2011. Privacy Preservation in the Dissemination of Location Data. *SIGKDD Explor. NewsL*. 13, 1 (2011), 6–18. <https://doi.org/10.1145/2031331.2031334>
- [46] U.S. General Services Administration. 2017. open data in USA. (2017). <https://www.data.gov/>
- [47] Marius Wernke, Pavel Skvortsov, Frank Dürr, and Kurt Rothermel. 2014. A classification of location privacy attacks and approaches. *Personal and Ubiquitous Computing* 18, 1 (2014), 163–175. <https://doi.org/10.1007/s00779-012-0633-z>
- [48] Yu Zheng, Xing Xie, and Wei-Ying Ma. 2010. GeoLife: A Collaborative Social Networking Service among User, location and trajectory. *IEEE Data(base) Engineering Bulletin* (2010).
- [49] Changqing Zhou, Dan Frankowski, Pamela Ludford, Shashi Shekhar, and Loren Terveen. 2004. Discovering Personal Gazetteers: An Interactive Clustering Approach. In *Proceedings of the 12th Annual ACM International Workshop on Geographic Information Systems (GIS '04)*. ACM, New York, NY, USA, 266–273. <https://doi.org/10.1145/1032222.1032261>