

Article

## Entropy-Based Privacy against Profiling of User Mobility

Alicia Rodriguez-Carrion <sup>1\*</sup>, David Rebollo-Monedero <sup>2</sup>, Jordi Forné <sup>2</sup>, Celeste Campo <sup>1</sup>, Carlos Garcia-Rubio <sup>1</sup>, Javier Parra-Arnau <sup>2</sup> and Sajal K. Das <sup>3</sup>

<sup>1</sup> Department of Telematic Engineering, University Carlos III of Madrid, Avda. Universidad 30, E-28911 Leganés, Madrid, Spain

<sup>2</sup> Department of Telematics Engineering, Universitat Politècnica de Catalunya (UPC), Campus Nord, C. Jordi Girona 1-3, 08034 Barcelona, Spain

<sup>3</sup> Computer Science Department, Missouri University of Science and Technology, 325B Computer Science Building, Rolla, MO 65409-0350, USA

\* Author to whom correspondence should be addressed; E-mail: [arcarrio@it.uc3m.es](mailto:arcarrio@it.uc3m.es); Phone: +34 91-624-6234.

Version May 21, 2015 submitted to *Entropy*. Typeset by *LaTeX* using class file *mdpi.cls*

---

1     **Abstract:** Location-based services (LBSs) flood the mobile phones nowadays, but their  
2     use pose an evident privacy risk. The locations accompanying the LBS queries can be  
3     exploited by the LBS provider to build the user profile of visited locations, which might  
4     disclose sensitive data such as work or home locations. The classic concept of entropy is  
5     widely used to evaluate privacy in these scenarios where the information is represented as a  
6     sequence of independent samples of categorized data. However, since the LBS queries might  
7     be sent very frequently, locations profiles can be improved by adding temporal dependencies,  
8     thus becoming mobility profiles, where location samples are not independent anymore and  
9     might disclose the user's mobility patterns. Since time dimension is factored in, the classic  
10    entropy concept falls short to evaluate the real privacy level, which depends also on the time  
11    component. Therefore, we propose to extend the entropy-based privacy metric to the use of  
12    entropy rate to evaluate mobility profiles. Then, two perturbative mechanisms are considered  
13    to preserve locations and mobility profiles under gradual utility constraints. We further use  
14    the proposed privacy metric and compare it to classic ones to evaluate both synthetic and  
15    real mobility profiles when the perturbative methods proposed are applied. The results prove  
16    the usefulness of the proposed metric for mobility profiles, and the need for tailoring the  
17    perturbative methods to the features of mobility profiles in order to improve privacy without  
18    completely losing utility.

**Keywords:** Location-Based Services (LBSs); entropy; privacy; perturbative methods; location history.

---

## 1. Introduction

Recent years have witnessed the growth of a rich variety of information and communication technologies. As a result, users enjoy applications striving to tailor information-exchange functionality to their specific interests. Examples of these applications comprise *location-based services (LBSs)*, personalized Web search, and multimedia recommendation systems. In the specific example of LBSs, they open up enormous business opportunities, encompassing intelligent personal assistants, concierge and emergency services, entertainment, and advertising among others. These services are also the result of recent advances in positioning techniques, such as *Global Positioning System (GPS)*, location-based social networks (where users tag significant places), or location techniques based on cellular or WiFi networks (where user's locations is inferred by associating the signal strength received by her mobile phone with information about the closest transmission points and their position).

Many of these location services build upon, or lend themselves to, the creation of user profiles. However, user profiles by themselves, but specially when combined across several information services, pose evident privacy and security risks. On the other hand, it is precisely the availability to a system of such sensitive information what enables such intelligent functionality. Therefore, the need for preserving privacy without compromising the utility of the information emerges naturally. Hard privacy [1] is one of the existing privacy-enhancing technologies (PETs) [2] that consists in the preservation of the privacy by the user itself by minimizing, obfuscating or perturbing the information released, without the requirement of trusted intermediaries. In principle, by perturbing the confidential data prior to its disclosure, users attain a certain degree of privacy, at the expense of degrading the system performance (or utility). The existence of this inherent compromise is a strong motivation to develop quantifiable metrics of privacy and utility, and to design practical privacy-enhancing, data-perturbative mechanisms achieving serviceable points of operation in this privacy-utility trade-off.

This work focuses on the specific scenario of the privacy risk associated to the profiling of user mobility arising from the use of LBSs and user locations based on the identifier of the cell to which the user's phone is attached to (as opposed to the numerous privacy analysis on GPS-based location). When using LBSs in mobile phones—such as weather, traffic, or news widgets, to name a few of the most basic services that depend on the user location—the user's phone sends, quite frequently, a service request together with the user location, aiming to obtain the most up to date information. For this kind of services, it is sufficient to know a coarse precision location, thus cell-based location being suitable whilst diminishing the battery consumption with respect to the GPS option. The LBS provider may, then, collect or disclose to third parties sensible data related to the locations visited by the user. In this work, we distinguish two types of profile that can be built from the collection of locations sent to the LBS provider. We define the first one as **location profile**, and consists of the set of locations visited

55 by the user and the visit frequency to each one. This profile may disclose implicit information related  
56 to the user: her home and work locations; if she has children (the number of visits to a kindergarten or  
57 school is high); if she may suffer from some chronic disease (the number of visits to a hospital is high);  
58 if she travels much (there are visits to locations located in many different countries), among others.  
59 In these cases, an attacker aims at obtaining the most accurate probability distribution of the visits to  
60 each location. Then, it would be easier for the attacker to derive the implicit information enclosed in  
61 the location profile if some few locations concentrate many more visits, i.e., if the location profile is as  
62 different as possible than an uniform distribution. There exist several metrics to measure privacy in this  
63 type of scenarios where a set of labeled data exposes the user profile. Some of them are based on the  
64 concept of entropy of a set of independent samples, but to the best of our knowledge it has never been  
65 applied to the specific case of sequences of cell-based locations.

66 Furthermore, we define a second type of profile that can be built by taking advantage of the frequent  
67 LBS requests mobile phones usually send to obtain the updated information related to their location.  
68 We denote it as **mobility profile**, and it is defined as the temporal sequence of locations visited by the  
69 user. Therefore, the stress in this profile lies on the correlations among the visited locations, instead of  
70 considering the locations as independent events. In this case, an attacker will aim at correctly predict  
71 the next location of the user, given her past history of locations. With this profile the adversary could  
72 derive more refined information due to the knowledge of temporal dependencies. An innocent example  
73 of personal mobility information disclosure might be the following one. If the untrusted LBS provider  
74 knows, by inspecting the mobility profile, that the user goes from home (first most visited location) to  
75 work (second most visited location) and then to a third location near a supermarket, the provider might  
76 infer that the user regularly buy products at that supermarket. Therefore, the LBS provider might leak  
77 this data to other related services which can start sending advertisement or offers of different shops  
78 offering the same products right before the user goes to her usual supermarket. This behavior, which  
79 might result very effective for advertisement, is only possible when adding the temporal component to  
80 the locations profile to transform it into a mobility profile. The problem arising in this situation is that  
81 not only the set of visited locations and their visit frequency are target of a privacy attack, but also  
82 the correlations among the visits to those locations constitute a privacy threat. As demonstrated in [3],  
83 the correlations among location samples enclose a great deal of information when aiming at predicting  
84 the next location of a user. Since this is the target of an adversary, we need to measure privacy taking  
85 into account such correlations. However, the classical concept of entropy used for the location profiles  
86 does not work on memory processes, because it is only applicable for sequences of independent samples.  
87 Therefore, applying privacy metrics based on entropy to a mobility profile does not reflect the real privacy  
88 level, since the temporal correlations among locations visits, which represent the main component of a  
89 mobility profile, remain ignored. To the best of our knowledge, there exist no privacy metrics addressing  
90 the extension from location profiles to mobility profiles. For this reason, we propose to compare the use  
91 of classical entropy with respect to an extension of this concept for processes with memory: the entropy  
92 rate. This general goal leads to the contributions stated in the following section.

### 93 *1.1. Contribution and Organization*

94 The main contributions of this work are the following ones:

- 95 • Jaynes' rationale on maximum entropy methods [4,5] enables us [6] to measure the privacy of  
96 confidential data modeled by a probability distribution by means of its Shannon entropy. In  
97 particular, one may measure the anonymity of a user's behavioral profile as the entropy of a relative  
98 histogram of online activity along predefined categories of interest. Inspired by this principle, we  
99 propose the use of Shannon's entropy to measure the privacy of a sequence of places of interest  
100 visited by a user (which we will refer to as the user's locations profile), with the caveat that this  
101 may only be appropriate for a series of statistically independent, identically distributed outcomes.
- 102 • Taking a step further, we tackle the case in which a sequence of location data is more adequately  
103 modeled as a stochastic process with memory, representing the (entire or recent) history of a  
104 user moving across predefined, discretized locations. We propose extending the more traditional  
105 measure of privacy by means of Shannon's entropy, to the more general information-theoretic  
106 quantity known as entropy rate, which quantifies the amount of uncertainty contained in a process  
107 with memory. In other words, we put forth the notion of entropy rates as the natural extension of  
108 Jaynes' rationale from independent outcomes to time series. Concordantly, we propose entropy  
109 rate as a novel, more adequate measurement of privacy of what we will call user mobility profiles:  
110 profiles capturing sequential behavior in which current activity is statistically dependent of the  
111 past, as it is commonly the case for location data.
- 112 • The extension from location to mobility profiles requires a reconsideration of the privacy  
113 preserving mechanisms. We propose two simple perturbative methods, previously used for Web  
114 search applications, looking for their suitability in these two profiling scenarios.
- 115 • Finally, we compare the results of calculating the privacy metrics proposed for different theoretical  
116 processes of increasing memory, to finally analyze a real location and mobility profile, made up  
117 of cell-based locations, which shows the usefulness of the proposed privacy metric. The work  
118 ends up with a discussion on different aspects impacting the privacy level obtained and further  
119 considerations to improve it in mobility profiling scenarios.

120 The remainder of the paper is organized as follows. Section 2 presents a detailed study of the state  
121 of the art in the two main topics covered in this work: hard-privacy, data-perturbative technologies;  
122 and privacy metrics for data perturbation against user profiling. Section 3 states the problem, taking  
123 care of the application scenario, the specific privacy model and metrics considered, as well as the  
124 perturbative methods to be used. Section 4 exposes the formal analysis of the problem at hand. In  
125 Section 5 the experimental data and results are described, together with a discussion on the findings and  
126 limitations found. Finally, Section 6 gathers the main conclusions along with some future lines, and the  
127 Appendixes A to C include the proofs to the mathematical expressions derived in Section 4.

## 128 2. Related work

129 As exposed in [7], the evolution of LBSs and the associated location techniques led to a privacy  
130 degradation. Anonymous location traces can be identified by correlation with publicly available

131 databases, thus increasing the possibility of disclosing sensitive data, such as home and work  
132 locations [8] or specific points of interest of the user [9]. Therefore, users are exposed to different  
133 kinds of attacks (e.g., tracking, localization or meeting attacks, among others [10]) with the available  
134 information collected by LBS providers. For this reason, privacy enhancement is key in order to tackle  
135 the increasing new threats that arise from the evolution of LBSs.

136 The following is a brief overview of the state of the art on privacy-enhancing technologies and privacy  
137 metrics related to LBSs and profiling of user mobility.

### 138 2.1. Privacy-Enhancing Technologies for LBSs

139 Many different privacy-enhancing techniques focused on LBSs and location profiling can be found  
140 in the literature. The Statistical Disclosure Control (SDC) community proposed many of them, aiming  
141 to prevent the disclosure of the contribution of specific individuals by inspecting published statistical  
142 information. *k-anonymity* [11,12] is one of the proposed techniques. A specific piece of data on  
143 a particular group of individuals is said to satisfy the *k-anonymity* requirement if the origin of any  
144 of its components cannot be ascertained, beyond a subgroup of at least  $k$  individuals. The concept  
145 of *k-anonymity* is a widely popular privacy criterion, partly due to its mathematical tractability.  
146 However, this tractability comes at the cost of important limitations, which have motivated a number  
147 of refinements [13–15].

148 In the context of statistical databases appears also the concept of **differential privacy** [16–18]. The  
149 idea behind this approach is to guarantee that, after adding random noise to a query, if it is executed  
150 on two databases that only differ on one individual, the same answer must be generated with similar  
151 probabilities in both databases. Differential privacy is used for LBSs when aggregate location data is  
152 published. However, our scenario is that of a single user sending requests to an LBS provider, which is  
153 a slightly different case. In order to cope with this difference, the concept of **geo-indistinguishability**  
154 emerged recently [19,20]. It is a variant of differential privacy for the specific case of LBSs based on the  
155 principle that, the closer two locations are, the more indistinguishable they should be. In other words,  
156 given two close locations, they should generate the same reported location to the LBS provider with  
157 similar probabilities.

158 Other widely used alternatives, known as user-centric approaches, rely on perturbation of the location  
159 information and **user collaboration**. In this last context, the authors in [21] propose the collaboration  
160 of the users to exchange context information among the interested user and another one who already has  
161 that piece of data. This way, many interactions with the LBS provider disappear, thus increasing the  
162 location privacy by avoiding as many requests (with the user's location attached to it) to the provider as  
163 possible. On the other hand, users interactions poses in some cases additional privacy risks. That is the  
164 case of the effect of co-location in social networks, as demonstrated in [22]. In these situations, even  
165 if the user does not disclose her location, she might reveal her friendship and current co-location with  
166 a user who does disclose her location. The authors then quantify the impact of these co-location data,  
167 deriving an inference algorithm.

168 Regarding the use of **location perturbation techniques**, we already introduced the concept of hard  
169 privacy [1,23], in its fundamental form of data perturbation carried out locally prior to its disclosure

170 (sometimes referred to as obfuscation), without the requirement of any trusted external party, but  
171 inducing a compromise between the privacy attained, and the degradation of the utility of the data  
172 disclosed for the intended purposes of an information service. A wide variety of perturbation methods  
173 for LBSs has been proposed [24]. We cannot but briefly touch upon a few recent ones. In [25], locations  
174 and adjacency between them are modelled by means of the vertices and edges of a graph, assumed to  
175 be known by users and providers, rather than coordinates in a Cartesian plane or on a spherical surface.  
176 Users provide imprecise locations by sending sets of vertices containing the vertex representing the actual  
177 user location. Alternatively, [26] proposes sending circular areas of variable center and radius in lieu of  
178 actual coordinates. Finally, we sketch the idea behind [27]. First, users supply a perturbed location,  
179 which the LBS provider uses to compose replies sorted by decreasing proximity. The user may stop  
180 requesting replies when geometric considerations guarantee that the reply closest to the undisclosed exact  
181 location has already been supplied. Besides these approaches, a number of hard-privacy mechanisms  
182 relying on data perturbation have been formulated in an application context wider than LBSs, primarily  
183 including online search and resource tagging in the semantic Web. Indeed, an interesting approach to  
184 provide a distorted version of a user's profile of interests consists in *query forgery*. The underlying  
185 principle is to accompany original queries or query keywords with bogus ones, in order to preserve user  
186 privacy to a certain extent. The associated cost relates to traffic and processing overhead, but on the  
187 other hand, the user does not need to trust the service provider nor the network. Building on this simple  
188 principle, several protocols, mainly heuristic, have been proposed and implemented, with various degrees  
189 of sophistication [28–30]. A theoretical study of how to optimise the introduction of bogus queries from  
190 an information-theoretic perspective, for a fixed constraint on the traffic overhead, appears in [31]. The  
191 perturbation of user profiles for privacy preservation may be carried out not only by means of insertion of  
192 bogus activity, but also by suppression [32]. These approaches constitute the basis of the present work.

193 Finally, going a step further by preserving not only privacy related to locations understood as a set  
194 of independent samples, but also the **correlations among locations**, the most recent works on location  
195 privacy like [33] take into account the sequential correlation between locations, aiming at protecting the  
196 present, past, and future locations, as well as the transitions between locations. The authors tackle the  
197 problem as a Bayesian Stackelberg problem, and use the attacker's estimation error as privacy metric.  
198 This problem is similar to our scenario, since the preserving the privacy of the correlations among  
199 locations is our main concern. However, we tackle the problem with a different approach, using entropy  
200 rate definition. On the other hand, whilst we do not propose any Location Privacy Preserving Mechanism  
201 (LPPM) (beyond a couple of naive approaches to demonstrate the usefulness of the proposed privacy  
202 metric), the authors of the mentioned work also defined a theoretical framework based on the Bayesian  
203 Stackelberg approach to preserve location privacy.

## 204 2.2. Privacy Metrics for Data Perturbation against User Profiling

205 Quantifiable measures of performance are essential to the evaluation of privacy-enhancing  
206 mechanisms relying on data perturbation, in terms of both the privacy attained and any degradation of  
207 utility. In a recent study on privacy metrics [34], it is shown that many of them may be understood from  
208 a unifying conceptual perspective that identifies the quantification of privacy with that of the error in the

209 estimation of sensitive data by a privacy adversary, i.e., privacy is construed as an attacker's estimation  
210 error.

211 Of particular significance is the quantity known as *Shannon's entropy* [35], a measure of the  
212 uncertainty of a random event, associated with a probability distribution across the set of possible  
213 outcomes.

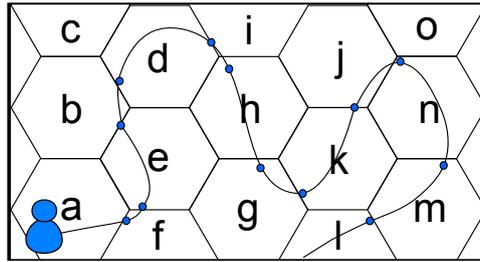
214 Some studies [36–41] propose the applicability of the concept of entropy as a measure of privacy, by  
215 proposing to measure the degree of anonymity observable by an attacker as the entropy of the probability  
216 distribution of possible senders of a given message in an anonymous-communication system. More  
217 recent works have taken initial steps in relating privacy to information-theoretic quantities [31,42,43].

218 A mathematically tractable model of *user profile* is a histogram of relative frequencies of visited  
219 locations, regarded as a probability distribution, on which we may compute information-theoretic  
220 quantities such as Shannon's entropy. In the focus of this paper, an intuitive justification in favor of  
221 entropy maximization is that it boils down to making the perturbed, observed user profile as uniform  
222 as possible, thereby hiding a user's particular bias towards certain visited places. A much richer  
223 argumentation stems from *Jaynes' rationale behind entropy maximization methods* [4,5], more generally  
224 understood under the perspective of the method of types and large deviation theory [35]. Under Jaynes'  
225 rationale on entropy maximization methods, the entropy of an apparent user profile, modeled by a relative  
226 frequency histogram, may be regarded as a measure of privacy, or more accurately, anonymity. The  
227 leading idea, proposed in [31,44], is that the method of types from information theory establishes an  
228 approximate monotonic relationship between the likelihood of a PMF in a stochastic system and its  
229 entropy. Loosely speaking and in our context, the higher the entropy of a profile, the more likely  
230 it is, and the more users behave according to it. This is of course in the absence of a probability  
231 distribution model for the probability mass functions, viewed abstractly as random variables themselves.  
232 Under this interpretation, entropy is a measure of anonymity, *not* in the sense that the user's identity  
233 remains unknown, but only in the sense that higher likelihood of an apparent profile, believed by an  
234 external observer to be the actual profile, makes that profile more common, hopefully helping the user  
235 go unnoticed.

### 236 3. Privacy-Enhanced Perturbation of Trajectories

237 In this section we describe the scenario where the privacy enhancement techniques will be applied, as  
238 well as the theoretical foundation of such techniques. First, we describe how mobility is represented in  
239 our scenario, highlighting the difference with respect to the GPS-based mobility representation. Next,  
240 once we have defined the data to protect and their representation, we need to know which mechanisms  
241 to use in order to enhance the user's privacy. But in order to evaluate the mechanisms, we first need to  
242 define how to measure the privacy level attained. This topic will lead to a discussion on how a concept  
243 such as entropy is a good privacy measure, and how to extend it to the domain of time series through the  
244 use of entropy rates. Finally, after defining a quantitative measure of privacy, we propose a perturbative  
245 method to enhance user's mobility data privacy under certain utility constraints.

#### 246 3.1. User Mobility Profiling and Adversary Model



**Figure 1.** Movement scenario divided in regions, and the trajectory followed by the user.

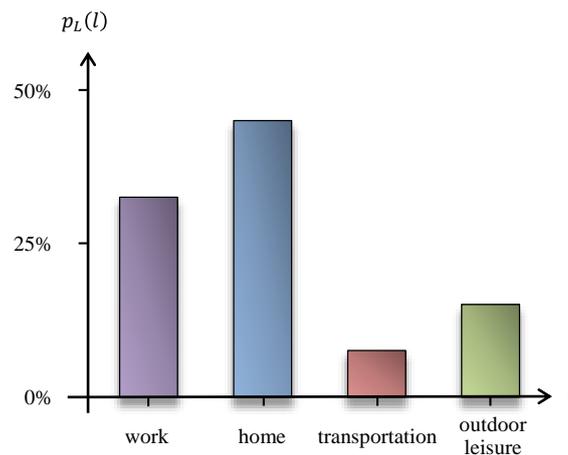
247 Users of an LBS disclose *trajectories*, i.e., sequences of positions, to a service provider. With small  
 248 loss of generality for the purposes of user profiling on the basis of behavior, we assume that those  
 249 positions are not treated in the form of space coordinates, but categorized into a predefined, finite set of  
 250 labeled symbolic locations. The movement scenario is divided into different regions, each one tagged  
 251 with a unique identifier. The user moves across this scenario, and each time she enters a different  
 252 region, the identifier corresponding to that region is recorded to what is known as location history  
 253 or trace,  $L$ . This kind of representation allows to record sequences such as locations represented by  
 254 GSM/UMTS cells, WiFi coverage areas or sequences of concrete places (office, home, market, gym. . .).  
 255 This assumption will enable us to model trajectories as random processes with samples distributed in  
 256 a finite alphabet. In Figure 1 we can see the track of a user that corresponds to the location history  
 257  $L = afdbdihgkijnml$ . Further, the data contained in a user's trace allow us to define two types of user  
 258 profile, the location profile and the mobility profile. In the following, we define and comment on these  
 259 two types of profiles, along with their corresponding adversary models.

260 • **Location profile.** This profile is defined as the probability distribution of the visits to each of the  
 261 locations in the set of visited locations of the user, i.e., the relative frequency of visits of the users  
 262 visited location set. This is analogous to the histogram of the relative frequency of the different  
 263 search categories, in the case of web search presented in [31,45]. This profile reveals information  
 264 related to different locations, independently of the rest of visited locations and correlations among  
 265 them. For instance, an attacker may be interested in knowing the probability distribution of the  
 266 visits in order to know several pieces of related data such as: home or work locations, which  
 267 are demonstrated to be very easy to derive [3,46] even when the attacker has access to just a  
 268 few LBS requests [47]; if the user travels to many different countries; if the user usually visits  
 269 (the relative visit frequency is high) some hospital, religious or political organization, children  
 270 school, sports center, among others. The attacker, say the LBS provider, or a third party to  
 271 whom the provider relinquishes the user location profile, might use this information to provide  
 272 personalized advertisement, or vary prices depending on the user's demand (e.g., if the frequency  
 273 of the cumulative visits to locations in a different country to the one with the highest number of  
 274 visits is high, it can be derived that the user travels frequently, thus she will be prone to book flights  
 275 at higher prices because traveling might be part of her work). A high number of visits to a hospital,  
 276 or a religious or political-related venue can have also impact when looking for jobs or insurances.

277 • **Definition.** Let  $L$  be a random variable (r.v.) representing the location of a given user, from  
 278 an alphabet of *predefined location categories*  $\mathcal{L}$ . The time of the location referred to is chosen  
 279 uniformly at random. We model the *location profile* of said user as the probability distribution

of  $L$ ; precisely, the probability mass function (PMF)  $p_L$  of the discrete r.v.  $L$ . Thus,  $p_L(l)$  is the probability that the user is at location  $l \in \mathcal{L}$  at any given time. In other words,  $p_L(l)$  represents the relative frequency with which the user visits this location.

- **Example.** Examples of location categories that may characterize the behavioral profile of a user include categories such as 'work', 'home', 'car', 'subway', 'restaurant', 'movie theatre', 'out of town'. These could be inferred from geographical locations with the help of an appropriate map. Fig. 2 depicts a simple example of location profile on a location alphabet with a few categories.



**Figure 2.** Example of location profile  $p_L$  as a PMF or histogram of relative frequencies on a simple location alphabet  $\mathcal{L} = \{\text{'work'}, \text{'home'}, \text{'transportation'}, \text{'outdoor leisure'}\}$ , inferred from geographical locations.

- **Adversary model for the location profile.** The adversary model is, in this case, estimating the visit probability distribution as accurately as possible, by inspecting the locations attached to the LBS requests. To this end, the adversary could utilize a maximum likelihood estimate of the distribution, directly as the histogram of relative frequencies, simply by counting observed locations, or any other well-known statistical techniques for estimation of probability distributions, such as additive or Laplace smoothing.

Our adversary model contemplates what the attacker is after when estimating those location profiles. According to [6], and in line with the technical literature of profiling [48,49], we assume the attacker's ultimate objective behind profiling is to target users who deviate significantly from the typical location profiles. This is known as *individuation*, meaning that the adversary aims at discriminating a given user from the whole population of users, or said otherwise, wishes to learn what distinguishes that user from the other users.

We would like to remark that our interpretation of Shannon's entropy as a measure of profile likelihood is clearly consistent with the assumptions made about the adversary, and in particular with its objective for constructing location profiles. Specifically, the higher the Shannon entropy of a location profile, the larger the number of users sharing this location pattern, and therefore the less interesting is the profile to an attacker assumed to target peculiar users. We hasten to stress that the Shannon entropy is, accordingly, a measure of *anonymity* rather than privacy, in the sense that the

obfuscated information is the uniqueness of the profile behind such location patterns, rather than the actual profile itself.

Another interpretation of Shannon's entropy as an anonymity metric stems from the intuitive observation that the higher the entropy of the distribution, informally speaking, the flatter the distribution, the less information the attacker could derive about predictable locations. In other words, if all the locations have the same visit frequency, the attacker can know the visited locations, but not which of them are more important.

- Mobility profile.** This profile is defined as the joint probability of visited locations over time, or equivalently, as the sequence of conditional probabilities of the current location, given the past history of locations. In this case, locations are not considered independently as in the user's location profile, but the most important component is the correlation among different locations, i.e., the short and long-range temporal dependencies among them. In this case, an attacker will aim at predicting the next location the user will visit, given the past history. The predictions about future locations provide a further refinement for advertisement purposes: the advertiser knows not only which product might be most interesting for the user regarding her visited locations, but also when to offer it for maximizing the impact of the add. For instance, suggesting some entertainment activity might be more effective if by the mobility profile the attacker finds the user did not go from home to work, as usual, which might indicate a weekend or holiday. The adversary goal is then to be able to predict as accurately as possible the next location of the user, given her past mobility history. There exists many prediction algorithms that can be used to do so [50], and their success depends on the predictability of the mobility history. As demonstrated in [3], the temporal dependencies among the locations visited by the user enclose information that noticeably increases the predictability of the mobility. In that study, the authors define the concept of predictability, closely linked to the entropy rate of the sequence, and demonstrate that it constitutes an upper bound on how much of the time we could correctly predict the next location of the user, given her past mobility history. After analyzing real mobility histories of thousands of users, they conclude that we could correctly predict the next location of a user 93% of the time in average.

- Definition.** More precisely, for each user we define a stochastic process  $(L_t)_{t=1,2,\dots}$  representing the sequence of categorized locations over discretized time instants  $t = 1, 2, \dots$ . Concordantly, the corresponding location  $L_t$  at time  $t$  is a discrete r.v. on the alphabet of predefined location categories  $\mathcal{L}$  introduced earlier. We define the *mobility profile* of the user as the joint probability distribution of locations over time,

$$p_{L_1 L_2, \dots, L_{t-1}, L_t, L_{t+1}, \dots}(l_1 l_2, \dots, l_{t-1}, l_t, l_{t+1}, \dots),$$

which may be equivalently expressed, by the chain rule of probabilities, as the sequence of conditional PMFs of the current location  $L_t$  given the past location history  $L_{t-1}, L_{t-2}, \dots$ , i.e.,

$$p_{L_t | L_{t-1}, L_{t-2}, \dots}(l_t | l_{t-1}, l_{t-2}, \dots).$$

**Table 1.** Main conceptual highlights of the adversary model assumed in this work.

<i>Who can be the privacy attacker?</i>	Any entity able to <i>profile</i> users based on their location and mobility patterns is taken into account. This includes the LBS provider and any entity capable of eavesdropping users' location data, e.g., Internet service providers, proxies, users of the same local area network and so on. Further, we also contemplate any other entity which can collect publicly available users' data. This might be case of an attacker crawling the location data of Twitter users <sup>1</sup> .
<i>How does the attacker model location profiles?</i>	The <i>location profile</i> of a user is modeled as the probability distribution of their locations within an alphabet of predefined location categories. Conceptually, a location profile is a histogram of relative frequencies of user location data across those categories.
<i>How does the attacker represent mobility profiles?</i>	The <i>mobility profile</i> of a user is modeled as the joint probability of visited locations over time. This model is equivalent to the sequence of conditional probabilities of the current location, given the past history of locations.
<i>What is the attacker after when profiling users?</i>	We consider the attacker wishes to <i>individuate</i> users on the basis of their location and mobility patterns. In other words, the adversary aims at finding users who deviate notably from the average and common profile.

Discretized times could be defined in terms of fixed time intervals, such as hours, or fractions thereof, or more simply but less informatively, as times relative to a change in activity of the user, so that the actual logged data is the order of the given locations in time, but not their duration.

- **Example.** Following up with the simple example of Fig. 2, with location alphabet  $\mathcal{L} = \{\text{'work'}, \text{'home'}, \text{'transportation'}, \text{'outdoor leisure'}\}$ , the mobility profile now incorporates time information, in the form of fixed time intervals, say 15 min. In this manner, one could record the average time spent at work, at home, on the road, or at various types of outdoor leisure activities, and also mobility patterns involving said locations. With a more detailed location alphabet, one may detect that a user predictably goes to work directly from home, or to the movies or to a restaurant after work on a given day of the week.
- **Adversary model for the mobility profile.** This datum gives us an idea of what an adversary could know about the future locations of a user when the mobility profile of this user is available, and raises the concern that motivates our contribution. Since predictability is directly linked with the *entropy rate* of the mobility profile as a stochastic process, rather than the entropy, (the higher the entropy rate, the lower the predictability, as shown in [3]), we could use this information theory concept in order to quantify the privacy of the user mobility profile in such a way that the less predictable a user is (the higher her entropy rate is), the higher her mobility profile privacy will be.

Analogously to the adversary model for location profiles, we also incorporate here the objective behind such profiling and assume an attacker that strives to find users with atypical mobility profiles. Jaynes' rationale behind the entropy-maximization method allows us to regard the entropy rate (formally presented next) as an anonymity metric that is consistent with this objective.

### 3.2. Privacy Model and Entropy Rate as a Privacy Measure

This work considers an abstract privacy model in which individuals send pieces of confidential data, related to each other in a temporal sequence, to an untrusted recipient. This intended recipient of the

358 data is not fully trusted. In fact, it is regarded as a privacy adversary capable of constructing a profile  
359 of sensitive user interests on the basis of the observed activity, or prone to leaking such observations to  
360 an external party who might carry out the profiling. Disclosure of confidential data to such untrusted  
361 recipient poses a privacy risk. However, it is precisely the submission of detailed data on preferences  
362 and activity which enables the desired, intelligent functioning of the underlying information system.  
363 Although this abstraction is readily applicable to a wide variety of information systems, we focus our  
364 exposition on the important example of LBSs.

365 We have mentioned in our review of the state of the art, Section 2.2, that the anonymity of a  
366 profile can be quantified in terms of the Shannon entropy of the probability distribution representing  
367 the histogram of relative frequencies profiling user behavior [6]. The cited work argues in favor of the  
368 use of this information-theoretic measure capitalizing on Jaynes' rationale for entropy maximization  
369 methods. Roughly speaking, Jaynes' argument boils down to postulate that high-entropy is more  
370 common than low-entropy. In the context of privacy, high-entropy profiles are more frequent and thus  
371 more anonymous.

372 More sophisticated user profiling may be carried out if the privacy adversary exploits the statistical  
373 dependence among location samples over time, in order to infer temporal behavioral patterns. This  
374 responds to the observation that the disclosure of a sequence of user locations poses a clear privacy risk,  
375 especially when these locations are viewed in conjunction and time is factored in. Examples include  
376 answers to questions such as, where does a user commonly go after work, before heading back home?  
377 On a typical weekend, what is the user's preferred activity after leaving their house? What route does  
378 the user usually follow to get to work or back home?

379 The natural extension of the measurement of anonymity by means of entropy to the case at hand,  
380 namely random processes with memory, is *entropy rate*, formally defined in Section 4.1. Because  
381 the definition of entropy rate is approximated by the entropy of a large block of consecutive samples  
382 (normalized by the number of samples), the very same argument in favor of entropy can be extended to  
383 entropy rate, the latter more suitable to user profiling in terms of trajectory patterns rather than individual  
384 locations.

385 We should remark that entropy has been often proposed as a privacy metric, on the intuitive basis  
386 that it constitutes a measure of uncertainty, even though formally speaking there exist many other  
387 such measures, R enyi entropy or variance, to name a couple. Even though we acknowledge the  
388 appropriateness of this intuition, we more formally resort to Jaynes' rationale on maximum entropy  
389 methods to argue, as in [6], that more entropic user profiles are also more common, and thus less  
390 idiosyncratic or characteristic of the specific habits of a particular individual. Consequently, the privacy  
391 metric proposed here, namely the entropy rate of the stochastic process modeling the sequence of  
392 discretized locations, represents a quantifiable measurement of the anonymity of the user, in the sense of  
393 commonness, rather than of the confidentiality of the data at hand.

### 394 3.3. *Perturbative Mechanisms*

395 Following the reasons stated in the introduction, particularly motivated by the advantages of  
396 hard privacy against the reliance on trusted intermediaries, we shall investigate theoretically and

397 experimentally two data-perturbative strategies prior to the disclosure of trajectories, in order to trade-off  
 398 usability for privacy. In the first strategy, referred to as *uniform replacement* from now on, with certain  
 399 probability, samples are replaced with values drawn according to a uniform distribution over the alphabet  
 400 of possible categorized locations. In the second mechanism, which will be called *improved replacement*  
 401 the same fraction of samples are replaced, although a more sophisticated policy is employed. Precisely,  
 402 the replacing samples are drawn from the distribution obtained from the solution to the problem for  
 403 optimized query forgery developed in [31]. We should point out that because the optimization carried  
 404 out was originally intended for memoryless processes and anonymity was measured by means of entropy  
 405 instead of entropy rate, the aforementioned improved solution need not be optimal whenever the privacy  
 406 attacker exploits existing statistical dependencies over time. Consequently, both mechanisms are merely  
 407 heuristics we choose to evaluate.

408 The probability of replacement is indicative of the degradation in data utility. As we will expose in  
 409 the next section, the theoretical analysis is equivalent for sample replacement and addition. In this last  
 410 case, the utility degradation is understood as an increase in the information sent to the LBS provider,  
 411 thus incrementing the energy consumption of the mobile device, and potentially, the economic cost of  
 412 data traffic. We consider here applications that can send location samples to the corresponding LBS  
 413 more frequently than in a normal situation (i.e., where no privacy-enhancing method is applied). That  
 414 allows to send fake locations together with the original ones without degrading the service provided,  
 415 only increasing the cost associated with a more intensive communication. From now on we will talk  
 416 about sample replacement, but keeping in mind that it could be extended to sample addition, by slightly  
 417 changing what we understand by utility in that case. Because sample values may occasionally be  
 418 replaced by themselves, especially if the number of location categories is small, counting the number  
 419 of effectively perturbed values is a more adequate measure of utility. While there is ample room for  
 420 the development of more sophisticated metrics of utility reflecting the quality of the LBS response, the  
 421 necessarily limited scope of this work prefers to cover the aspects of privacy and perturbation, as a first  
 422 insightful step towards the problem of privacy-enhanced perturbation of processes with memory.

#### 423 **4. Theoretical Analysis of Perturbative Methods and Entropy-based Privacy Metric**

##### 424 *4.1. Notation and Information-Theoretic Preliminaries*

425 Throughout the paper, we shall follow the convention of uppercase letters for *random variables* (r.v.'s),  
 426 and lowercase letters for particular values they take on. For simplicity, all r.v.'s in this analysis take on  
 427 values in a *finite alphabet*. *Probability mass functions* (PMF) are denoted by  $p$ , subindexed by the  
 428 corresponding name of the r.v. when not understood from the context. For instance, we may denote the  
 429 PMF of an r.v.  $X$  at  $x$  by  $p_X(x)$ , or simply by  $p(x)$ .

We review a few fundamental results from information theory. The reader may refer to [35] for  
 specific details and proofs. The Shannon *entropy* of a r.v.  $X$  with PMF  $p$  and finite alphabet  $\mathcal{X}$  is  
 written interchangeably as  $H(X)$  or  $H(p)$ . Recall that entropy is maximized for the uniform distribution,  
 and for this distribution only, and that the *maximum* attained is the logarithm of the cardinality of the  
 alphabet. Put mathematically,  $H(p) \leq \log |\mathcal{X}|$ , with equality if and only if  $p$  is the uniform distribution.

Throughout this work all logarithms are taken to base 2, and subsequently the entropy units are *bits*. Recall also that  $H(p)$  is a *strictly concave* function of  $p$ , in the sense that for any distributions  $p$  and  $p'$  over the same alphabet, and any  $\lambda \in [0, 1]$ ,

$$H((1 - \lambda)p + \lambda p') \geq (1 - \lambda) H(p) + \lambda H(p'),$$

430 with equality if and only if  $\lambda = 0$ ,  $\lambda = 1$ , or  $p = p'$ .

Let

$$(X_n)_{n \in \mathbb{Z}} = \dots, X_{-2}, X_{-1}, X_0, X_1, X_2, \dots$$

be a *stationary random process* with samples defined on a common alphabet  $\mathcal{X}$ . Stationarity implies that both the entropy sequences  $\frac{1}{n} H(X_1, \dots, X_n)$  and

$$H(X_n | X_1, \dots, X_{n-1}) = H(X_1 | X_2, \dots, X_n)$$

431 are non-increasing and have a common limit, called *entropy rate*, denoted here by  $H_R(X)$ . For  $n$  large,  
 432 either of these entropy quantities constitutes an arbitrarily accurate approximation to the entropy rate  
 433 of the process. We can compute these quantities by choosing an appropriate value of  $n$ , such that the  
 434 blocks capture the correlations of the process, and calculating  $p_X(X_n | X_1, \dots, X_{n-1})$  as the number of  
 435 blocks  $X_1, \dots, X_n$  normalized by the total number of blocks of length  $n$ , and then applying the previous  
 436 formula  $H_R(X) = - \sum p_X(X_n | X_1, \dots, X_{n-1}) \log p_X(X_n | X_1, \dots, X_{n-1})$ .

437 Stationarity also implies that the samples of the process are identically distributed according to a  
 438 common PMF. When in addition, they are statistically independent, the process, or the samples thereof,  
 439 are then called *independent, identically distributed* (i.i.d.). More colloquially, a process with independent  
 440 samples is called *memoryless* or without memory. For an i.i.d. process, entropy rate and the entropy of  
 441 individual samples coincide, that is,  $H_R(X) = H_R(X_n)$ . For a general stationary process  $H_R(X) \leq$   
 442  $H_R(X_n)$ , with equality if and only if the process is memoryless. The highest entropy rate is attained by  
 443 processes with independent, uniformly distributed samples, that is,  $H_R(X) \leq \log |\mathcal{X}|$ , with equality if  
 444 and only if the process is uniformly distributed and memoryless.

#### 445 4.2. Perturbative Mechanisms

446 Again, consider a *stationary random process*  $(X_n)_{n \in \mathbb{Z}}$  with samples distributed on a common *finite*  
 447 alphabet  $\mathcal{X}$ . We shall argue elsewhere that entropy rate is an appropriate *privacy measure*. We propose  
 448 two alternative privacy-enhancing *perturbative mechanisms*, in which individual samples of the random  
 449 process  $X_n$  are replaced with  $X'_n$ , with probability  $\rho$  and independently from each other, as follows.

- 450 • Uniform replacement:  $X'_n$  is drawn uniformly from  $\mathcal{X}$ .
- 451 • Improved replacement:  $X'_n$  is drawn according to the distribution obtained as the solution to the  
 452 maximum-entropy problem of [31].

453 Even though [31] was meant for sample addition rather than replacement, the mathematical formulation  
 454 turns out to be completely equivalent. However, we should hasten to point out that the optimality  
 455 guarantee of the cited work applies to entropies of individual samples, but *not* entropy rates in general

456 processes with memory. Consequently, the two alternative mechanisms described above are merely  
 457 heuristic in the context of this work. In both cases, the resulting *perturbed* process  $(X'_n)_{n \in \mathbb{Z}}$  is clearly  
 458 stationary.

459 We shall call  $\rho$  the *replacement rate*. Because sample values may be conceivably replaced with  
 460 themselves, we propose the following *utility measure*, which more accurately reflects the actual impact  
 461 of the perturbative mechanism. Precisely, we define the *perturbation rate*  $\delta = \mathbb{P}\{X_n \neq X'_n\}$ , constant  
 462 with  $n$  on account of the stationarity of the processes involved, and observe that  $\delta \leq \rho$ , as only replaced  
 463 samples may be effectively perturbed, that is, actually different.

464 Even in the heuristic called improved replacement, the samples to be replaced are chosen randomly  
 465 and replaced independently of their original value. A truly optimal strategy, however, should choose  
 466 which samples to replace, exploit the statistical model of the memory of the process, and be optimized  
 467 for  $\delta$  rather than  $\rho$  as a measure of utility. The scope of this work is limited to the analysis of the  
 468 heuristic mechanisms described, as a first step towards shedding some light on the problem of designing  
 469 perturbative strategies for processes with memory and with a truly optimal privacy-utility trade-off (or  
 470 privacy-cost qualitatively talking if we would consider sample addition).

#### 471 4.2.1. Uniform Replacement

472 We prove that uniform replacement on stationary processes with a strictly positive replacement  
 473 rate will always increase the entropy rate, unless the original process is uniformly distributed and  
 474 memoryless.

475 **Lemma 1.** *Let  $S$  and  $U$  be independent r.v.'s, the latter uniformly distributed on the alphabet of the*  
 476 *former. Let  $T$  be a third r.v., in general statistically dependent on  $S$ . Take  $S' = U$  with probability  $\rho$ ,*  
 477 *independently from  $S$  and  $T$ , and  $S' = S$  otherwise. Then,  $\mathbb{H}(S'|T) \geq \mathbb{H}(S|T)$ , with equality if and only*  
 478 *if either  $\rho = 0$ , or else  $S$  is uniform and independent of  $T$ . (Refer to Section A for the demonstration)*

**Theorem 2.** *Let  $X = (X_n)_{n \in \mathbb{Z}}$  be a stationary random process with samples distributed on a common*  
*finite alphabet  $\mathcal{X}$ . Although the process  $X$  itself need not be independent, each of its samples  $X_n$  is*  
*altered completely independently as follows. Each sample  $X_n$  is replaced by another r.v.  $U_n$ , uniformly*  
*drawn from the alphabet  $\mathcal{X}$ , with probability  $\rho$ , and left intact otherwise. Let  $X' = (X'_n)_{n \in \mathbb{Z}}$  be the*  
*resulting process, also stationary. Then, for any  $m \geq 0$ ,*

$$\mathbb{H}(X'_0|X'_{-1}, \dots, X'_{-m}) \geq \mathbb{H}(X_0|X_{-1}, \dots, X_{-m}),$$

479 *with equality if and only if either  $\rho = 0$ , or else  $X_0$  is uniform and independent of  $X_{-1}, \dots, X_{-m}$ .*  
 480 *The same inequality holds in the limit of  $m \rightarrow \infty$  yielding entropy rates, that is,  $\mathbb{H}(X') \geq \mathbb{H}(X)$ ,*  
 481 *with equality if and only if either  $\rho = 0$ , or else  $X$  is uniformly distributed and memoryless. (Refer to*  
 482 *Section B for the demonstration)*

#### 483 4.2.2. Uniform versus Improved Replacement

484 We show that in the case of memoryless processes not originally uniform, improved replacement  
 485 will require a lower replacement rate to achieve maximum entropy than that demanded by uniform

486 replacement. We shall also see that when the cardinality of the alphabet is large, the perturbation rate  
487 approaches the replacement rate.

In the perturbative mechanisms described earlier, define the *critical replacement rate*  $\rho_{\text{crit}}$  to be the replacement rate  $\rho$  required for the entropy rate  $H(X')$  of the perturbed process  $(X'_n)_{n \in \mathbb{Z}}$  to attain its maximum possible value  $\log |\mathcal{X}|$ , achievable only when  $X'$  becomes memoryless and uniformly distributed. Denote by  $\delta_{\text{crit}}$  the corresponding, *critical perturbation rate*. Write

$$p_{\max} = \max_{x \in \mathcal{X}} p(x) \geq \frac{1}{|\mathcal{X}|},$$

488 with equality if and only if  $X$  is uniformly distributed.

489 **Theorem 3.** *Assume the nontrivial case in which the original process  $X$  is not already independent,*  
490 *uniformly distributed.*

*In uniform replacement,*

$$\delta = \rho \left(1 - \frac{1}{|\mathcal{X}|}\right),$$

$$\rho_{\text{crit}} = 1,$$

$$\delta_{\text{crit}} = 1 - \frac{1}{|\mathcal{X}|}.$$

*In improved replacement, for any  $\rho \geq 1 - \frac{1}{|\mathcal{X}|p_{\max}}$ ,*

$$\delta = (1 - \rho) \sum_x p(x)^2 + \rho - \frac{1}{|\mathcal{X}|}.$$

*If the original process is i.i.d.,*

$$\rho_{\text{crit}} = 1 - \frac{1}{|\mathcal{X}|p_{\max}},$$

$$\delta_{\text{crit}} = 1 - \frac{1}{|\mathcal{X}|} - \frac{1}{|\mathcal{X}|p_{\max}} \left(1 - \sum_x p(x)^2\right).$$

*Otherwise, in the general case of processes with memory,*

$$\rho_{\text{crit}} = 1 \text{ and } \delta_{\text{crit}} = 1 - \frac{1}{|\mathcal{X}|}.$$

491 *(Refer to Section C for the demonstration)*

Recall [35] that the *Rényi entropy* of order  $\alpha$  of a discrete r.v.  $X$  with PMF  $p_X$  is defined as

$$H_\alpha(X) = \frac{1}{1 - \alpha} \log \sum_x p_X(x)^\alpha. \quad (1)$$

492 The value  $\sum_x p(x)^2 = \mathbb{E} p(X)$  in the theorem is directly related to the Rényi entropy of order 2 of  
493  $p(x)$ , called collision entropy. The sum of squared probabilities above is minimized for the uniform  
494 distribution, and maximized for a degenerate distribution, where the associated r.v. takes on a single  
495 value with probability 1.

496 Observe that in the case of uniform replacement, a large alphabet  $|\mathcal{X}|$  implies that the perturbation  
497 rate will approach the replacement rate, that is,  $\delta \simeq \rho$ , because the unlikelihood of replacing a sample  
498 by itself. In the case of improved replacement, the approximation requires not only  $|\mathcal{X}| \gg 1$ , but also  
499  $\sum_x p(x)^2 \ll 1$ , and only holds for sufficiently large  $\rho$ .

500 *4.3. Entropy Estimation*

501 As previously shown, entropy could be a suitable privacy metric, but we should pay attention to the  
 502 estimator used. Depending on the concrete application or data to focus on, entropy estimation might be  
 503 different. In the case of human mobility, the location traces (that lead to locations and mobility profiles)  
 504 have specific features to take into account when estimating their entropy: strong long-range time-space  
 505 dependencies, high probabilities of returning to some highly frequented locations [46], the high number  
 506 of different visited places (cardinality of the alphabet), among others.

507 Bearing these features in mind, we could come up with different entropy estimates, as described in [3],  
 508 each one of them accounting for different dependencies. As we shall see next, two of these estimates  
 509 will be the Hartley entropy and the Shannon entropy. Throughout this subsection, these entropies will  
 510 be denoted by  $H_0(X)$  and  $H_1(X)$  to emphasize their connection with the Rényi entropy, a family of  
 511 functionals widely used in information theory as a measure of uncertainty. Particularly, from (1) it is  
 512 straightforward to see that, when  $\alpha = 0$ , Rényi's entropy boils down to Hartley's. In the limit when  $\alpha$   
 513 approaches 1, this family of functionals reduces to Shannon's entropy.

- Hartley entropy,  $H_0(X)$ , is the maximum attainable entropy value. We should recall that entropy is maximized for the uniform distribution, and for this distribution only, and that the *maximum* attained is the logarithm of the cardinality of the alphabet. Put mathematically:

$$H_0(X) \leq \log |X| \quad (2)$$

514 with equality if and only if  $X$  is drawn from the uniform distribution. Applied to our case, it would  
 515 be calculated considering the probability mass function of the locations trace (since no temporal  
 516 dependencies are considered) to be a uniform distribution of  $\mathcal{X}$  different symbols (locations).  
 517 This entropy represents the highest possible uncertainty, as it does not take into account temporal  
 518 aspects nor the number of visits accumulated by each location.

- Shannon entropy,  $H_1(X)$ , is calculated as:

$$H_1(X) = - \sum_i p_X(x_i) \log p_X(x_i) \quad (3)$$

519 In our scenario,  $p_X(x_i)$  is the probability of visiting location  $x_i$ , which can be computed as  
 520  $p_X(x_i) = \frac{N_{x_i}}{N}$ , where  $N_{x_i}$  is the number of visits received by location  $x_i$ , and  $N$  is the total number  
 521 of visits (i.e. the length of the movement history). Shannon entropy considers the correlations in  
 522 the location visits frequencies, thus being lower (or equal if the probability to visit each location  
 523 is the same) than  $H_0(X)$ . Actually, this entropy would be lower than  $H_0(X)$  as less uniform is  
 524 such PMF (i.e., as some locations receive many more visits than other ones). Locations profiles  
 525 (because no temporal dependencies are considered) behave precisely like this: some locations  
 526 corresponding to home or work unite the majority of visits, whilst the rest of locations are much  
 527 less visited.

- Entropy rate,  $H_R(X)$ , comes to scene when dealing with stationary processes, as pointed out in Section 4.1. It takes into account temporal dependencies between samples of the mobility profile

(in this case we consider the mobility instead of the locations profile because the time dependencies must be considered). Since  $H_R(X)$  takes into account more correlations of the profile, the resulting value is lower than the previous ones (there is less uncertainty regarding the next symbol of the profile), or equal if there are no temporal dependencies.

Applied to our case, we have a finite number of samples of the profile. Therefore, in order to obtain a good estimate of  $H_R(X)$ , we should choose the optimal block size,  $n$ . This block size should be large enough so that the blocks include important long-term temporal dependencies among location samples. But, since the length of the process (i.e., mobility profile) is limited and the cardinality of the alphabet (i.e., number of different locations) is high, there are no many samples of long blocks. Thus, choosing a block size too large leads to a poor estimate of  $p_X(X_1, \dots, X_n)$ . In order to use an appropriate value of  $n$ , we could use a well known entropy rate estimator based on Lempel-Ziv compression algorithms [3,51]. This way, the estimate of the entropy rate can be calculated as:

$$H'_R = \frac{\ln N}{N \sum_i \Delta_i} \quad (4)$$

where  $\Delta_i$  is the shortest substring starting at position  $i$  which has not been seen before from position 1 to  $i - 1$ , and  $N$  being the number of samples of the profile.

## 5. Experimental Study: Results and Discussion

In the previous section we formulated the theoretical problem of privacy-enhancing in processes with and without memory and how we tackle it. In [31] the authors show some results when the mechanisms proposed are applied to web queries, memoryless process, and using a small number of categories. In this section we will see what happens when the scenario switches to LBSs, where the number of categories increases, the probability model underneath becomes more complex and time starts playing an essential role. First, we will show the privacy gain obtained after applying the privacy enhancing mechanisms to different processes, both synthetic and real, and finally we discuss the differences that using real location data brings to the generic problem.

### 5.1. Results

This section collects the results drawn from applying the perturbative mechanisms described in the previous section to two different datasets. On the one hand, we will use several symbol sequences generated from Markov processes and basic alphabets of 2 symbols. With these data we will check the performance of the perturbative methods in simple ideal conditions and observe the influence of an increase of the process memory. On the second hand, real traces, taken from the Reality Mining dataset [52] will be processed and compared with the results of Markov processes, since the real scenario can be considered as an extrapolation of simple Markov processes in terms of memory and cardinality of the alphabet. More precisely, the locations history considered collects the sequence of locations visited by a user (each location represented as the identifier of the cell the user's phone was attached to at each instant) during an academic year. It gathers more than 500 different cells (symbols) and more than 10,000 cell changes (profile samples or location history length).

557 In order to show the privacy enhancement evolution, each process is perturbed from 0% of replaced  
 558 samples (i.e., the original symbol sequence) to 100% of replacements (all samples are replaced), as  
 559 explained in [31]. For each process and percentage of replacements, 10 realizations are averaged. As a  
 560 general rule, when  $\rho = 0$ , we have the original process and therefore the original (and minimum) entropy  
 561 value. As  $\rho$  increases, the process starts to become a uniform distribution, which is reached when  $\rho$  is  
 562 maximum, i.e., when all samples are replaced by another one using the perturbative methods previously  
 563 described, and therefore for the maximum value of  $\rho$ , the entropy value should be equal to  $H_0$ . We should  
 564 recall that  $\rho$  is the percentage of replacements, but since sometimes the replaced sample is equal to the  
 565 original one, the real replacement rate is  $\delta$ .

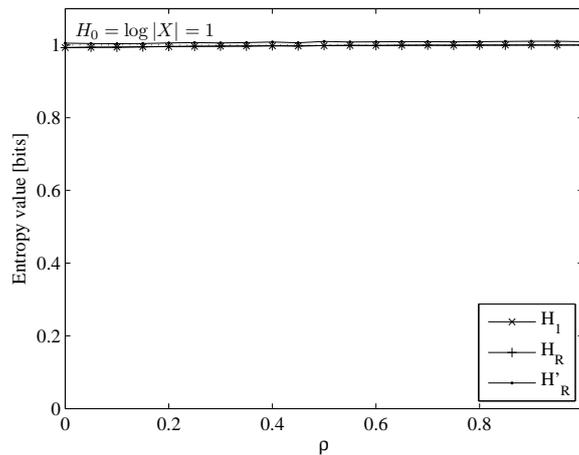
566 First, we will analyze how different entropy estimates work when applied to different kind of  
 567 processes. More precisely, we will study the influence of an increase of the process memory in the  
 568 entropy estimation as well as what happens when the process gets away from a uniform distribution,  
 569 both in terms of Shannon entropy and entropy rate. For this last case, we will compare two approaches:  
 570 the estimation by means of block entropies and the Lempel-Ziv-based estimate.

571 Figures 3 to 6 shows four different processes of 10000 samples in each of the plots:

- 572 • An almost uniform distribution, drawn from an order-1 Markov process with  $p(1|0) = 0.45$ ,  
 573  $p(0|1) = 0.55$ ,  $p(1) = 0.55$ ,  $H_0 = H_1 = H_R = 0.993$ . This is the base case.
- 574 • An i.i.d. (not uniform) process, drawn from an order-1 Markov process with  $p(1|0) = 0.8$ ,  
 575  $p(0|1) = 0.2$ ,  $p(1) = 0.8$ ,  $H_0 = 1$ ,  $H_1 = H_R = 0.772$ . Here we keep the process memoryless, and  
 576 change the probability distribution such that there is a bias towards one of the two symbols of the  
 577 alphabet.
- 578 • A Markov process with  $p(1|0) = 0.2$ ,  $p(0|1) = 0.05$ ,  $p(1) = 0.8$ ,  $H_0 = 1$ ,  $H_1 = 0.772$ ,  $H_R =$   
 579  $0.374$ . In this case we increase the memory of the process, keeping the cardinality and probability  
 580 distribution with respect to the second case.
- 581 • A real mobility trace taken from the dataset provided by the Reality Mining Project [52]. We can  
 582 only theoretically know  $H_0 = 8.765$  (drawn from the cardinality of the alphabet, i.e., the number of  
 583 different symbols representing the locations visited by the user), since the underlying probability  
 584 distribution is unknown. This means an increase both in the cardinality and the memory of the  
 585 process, due to the long range dependencies of human mobility.

586 For each figure (process) the entropy value evolution with respect to the replacement rates is plotted.  
 587 The samples are replaced using the uniform perturbative method, i.e., choosing the new sample from  
 588 the original alphabet of the sequence with the symbols uniformly distributed. Each process has been  
 589 generated 10 times, and the results shown here are the average value of the entropy calculated in each  
 590 repetition.

591 In the first case shown in Figure 3, the process without replacements is already uniform, therefore  
 592 there is no evolution in none of the entropy estimates. When the process is not uniform but still i.i.d.,  
 593 such as the one in Figure 4,  $H_1$  and  $H_R$  coincide, as there is no temporal information that can be captured  
 594 by  $H_R$  to lower the uncertainty, but they are lower than  $H_0$  and increase as the replacements turn the  
 595 process into a uniform one.



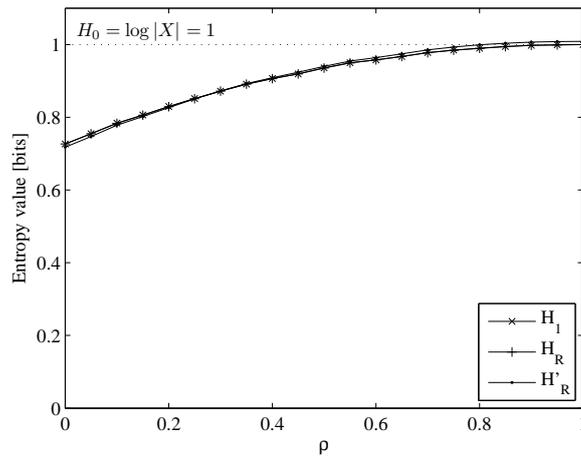
**Figure 3.** Different entropies for a process drawn from uniform distribution.

596 Figure 5 shows what happens when the process is not i.i.d. anymore. In this case,  $H_R$  is lower than  
 597  $H_I$  as it leverages the temporal information present now in the process to lower the uncertainty.

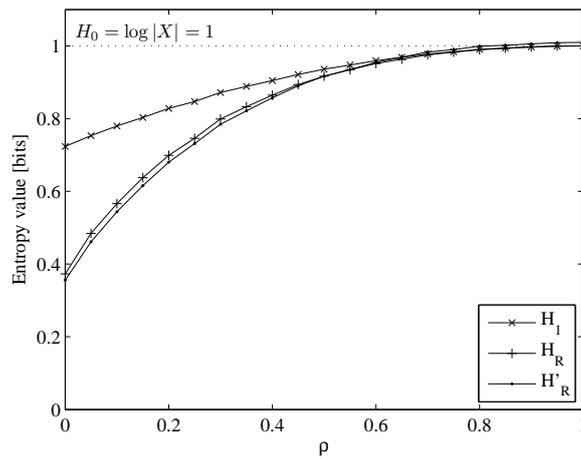
598 Finally, in Figure 6 we can see what happens when the number of different symbols (locations)  
 599 increases, as well as the memory of the process. In this case, we have 500 different symbols, what  
 600 leads to  $\frac{500!}{498!} = 249500$  possible blocks of 2 symbols to compute  $H_R$  using block entropies (the blocks  
 601 are of two symbols to compare with respect to the Markov processes). As the number of possible  
 602 blocks is so high and the number of samples is only of 10000, as the process becomes uniform, more  
 603 different blocks of two symbols come to scene. With this number of samples we do not have even one  
 604 occurrence of each different block, which probability would be  $p_X(X_1, \dots, X_n) = \frac{498!}{500!} = 4 * 10^{-6}$ .  
 605 Therefore, when computing  $H_R(X)$ , the values of the elements of the summation are very small, due  
 606 to the scarcity of occurrences of each possible block. This scarcity becomes more severe as the process  
 607 tends to uniformity. Thus,  $H_R(X)$  decays to near zero as the number of replaced samples increases, as  
 608 shown in the figure. As we previously explained in Section 4.3, this entropy estimation is biased by the  
 609 small number of samples available in the location history of the user (even when it comes from a year of  
 610 location tracking). This is the reason behind considering a different estimator like the Lempel-Ziv-based  
 611 one. Figure 6 shows how this estimator obtains more reasonable results. Both  $H_R(X)$  and  $H'_R(X)$  are  
 612 equal for the original sequence ( $\rho = 0$ ). However, in order to analyze the privacy improvement, we need  
 613 an estimate that works well for all the replacement rate span. We can also observe that, as the cardinality  
 614 of the alphabet is much higher, it is more difficult to choose the same sample as the original one in each  
 615 replacement, and therefore  $\delta$  is not bounded to 0.5 as in Markov processes of two symbols.

616 Now that we know how each entropy estimate works for different processes, it is time to apply such  
 617 estimates to the problem of privacy enhancement.

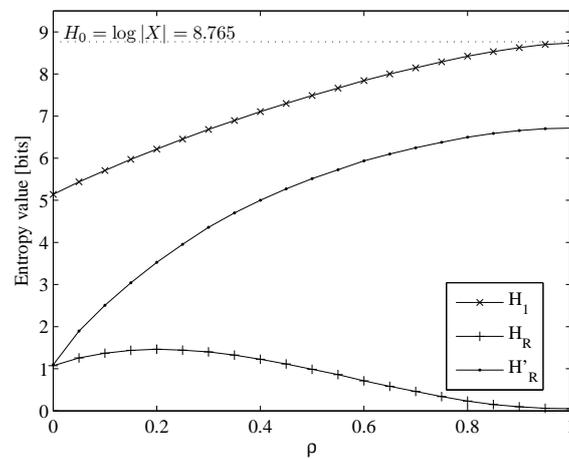
618 Figures 7 and 8 represent the privacy level obtained using the perturbative methods described  
 619 in Section 3.3, for the third Markov process described before (the one with memory) and for the  
 620 mobility history, respectively. In each figure, 4 plots can be distinguished: the privacy enhancement  
 621 in terms of entropy value, both for Shannon and entropy rate estimates, and for the two perturbative  
 622 methods considered. Measuring the privacy enhancement by means of two entropy estimates allows to



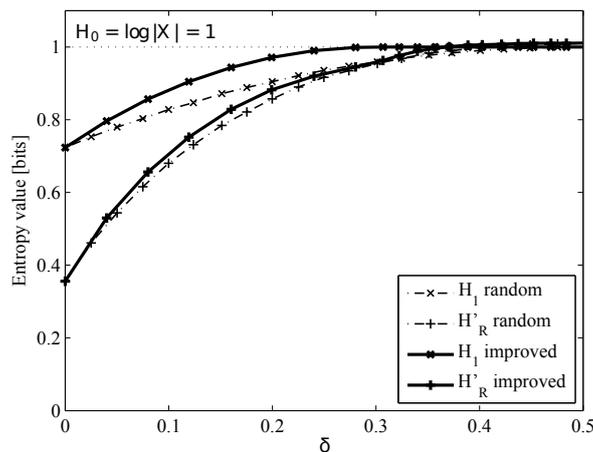
**Figure 4.** Different entropies for a process drawn from i.i.d. distribution.



**Figure 5.** Different entropies for a process drawn from Markov chain.



**Figure 6.** Different entropies for a process drawn from real mobility trace.



**Figure 7.** Comparison of perturbative methods for different privacy measures in a Markov process.

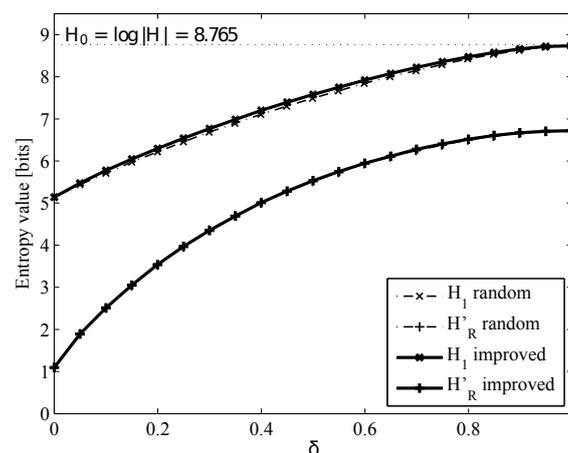
623 differentiate the results when only frequency-based information is considered (Shannon's entropy) from  
 624 the conclusions drawn from time-based data (entropy rate estimate).

625 For the case of the Markov process with memory in Figure 7, we can see that the privacy enhancement  
 626 is faster for the improved perturbative method, but only when no time-based information is considered.  
 627 Besides, it reaches the maximum privacy level (entropy value) when 35% of samples are replaced, value  
 628 that lowers up to 25% when the improved perturbative method is used and no temporal information is  
 629 considered.

630 When this same analysis is applied to the mobility trace, the results are quite different, as shown in  
 631 Figure 8. In this case, as the cardinality of the alphabet is so high, it requires a 100% of replacements  
 632 in order to obtain the highest privacy level, when measuring privacy as Shannon's entropy. Besides, the  
 633 maximum entropy is only achieved when no temporal correlations are considered. In order to get the  
 634 maximum value for sequence-based data, many more samples would be needed in order to have precise  
 635 entropy estimation. In this case it could be checked that the improved perturbation method does not  
 636 provide faster privacy enhancement for any case.

## 637 5.2. Discussion

638 In the previous figures it could be observed the great difference between theory, with simple Markov  
 639 processes, and real scenarios, such as users mobility profiles. But, where do these differences stem from?  
 640 Although the high cardinality of the alphabet and the complexity of the short and long term dependencies  
 641 of location histories play an important role, the probability distribution underneath the mobility trace is  
 642 also crucial. A great majority of visits are concentrated in two or three locations, corresponding to  
 643 home, work and the main points of interest of the user. Therefore the probability distribution is very  
 644 biased toward certain locations. The improved perturbative method is based on flattening the underlying  
 645 distribution with as few replacements as possible in order to get closer to a uniform distribution, and thus  
 646 to maximum entropy (i.e. privacy). When the number of categories is not very high and the probability  
 647 distribution is not very biased to certain few categories, it is easier to flatten it, as in the case of the



**Figure 8.** Comparison of perturbative methods for different privacy measures in a real mobility trace.

648 Markov processes shown. However, in order to flatten the mobility traces, we would have to compensate  
 649 the visits to 2 or 3 locations through the rest of the 500 different locations visited along the year. Although  
 650 there are more than 10000 samples, the cardinality is still very high, and would need many more samples  
 651 to be flattened. This issue is even more critical when considering not the distribution of the visits, but  
 652 the sequences of locations. If we consider short-term dependencies (short sequences) we are neglecting  
 653 important information, and even in this case the number of combinations is too high to compensate  
 654 the number of occurrences of the most repeated sequences. Considering long-term dependencies (long  
 655 sequences) leads to so many combinations that there are not enough samples to even calculate a good  
 656 entropy estimate, even worse if we try to flatten the block probability distribution.

657 The bias in the visits probability distribution carries an important consequence: for an attacker, it is  
 658 very easy to analyze a set of locations and determine where the main points of interest of the user are.  
 659 Therefore, these are very sensible data that must be masked. The bias can be leveraged in such a way  
 660 that, instead of trying to flatten all the distribution, we could focus on the set of the most visited locations  
 661 and just flatten their number of visits, leaving the least visited ones as they are. This way the uncertainty  
 662 of which of the most visited locations is home or work increases with few number of replacements. If  
 663 the number of replacements is not critical, or if we could fake the locations to be disclosed, the approach  
 664 could be to select some of the least visited locations and increase their number of visits to make it  
 665 comparable to the most visited places. However, as mentioned before, this strategy will require a great  
 666 number of replacements or additional fake locations. We should remark that adding fake locations incur  
 667 in a battery and data traffic increase, thus being utility-related factors to be taken into account when  
 668 deciding which perturbative approach to follow.

669 In the case of location sequences, i.e., focusing on preserving the privacy of the correlations among  
 670 locations, improving such privacy without compromising the data utility (or avoiding additional cost  
 671 when adding fake samples instead of replacing the original ones) is more complicated and depends  
 672 heavily on the application at hand. As we observe in the figures, in order to obtain high privacy levels,  
 673 the fraction of location samples to change grows fast. Furthermore, the replacements should be done  
 674 wisely. For example, let be a user walking in Madrid. If during the user's location sampling done by

675 the corresponding mobile application communicating with the LBS provider (done every few minutes)  
676 the system replaces a location in Madrid by another one in New York city (or just adds the location in  
677 New York city in the mobility profile), an attacker could easily detect that it is not possible for a user to  
678 make this large jump in such a short period of time. Therefore this replacement/addition might seem to  
679 theoretically improve greatly the privacy level (it is an unexpected movement, thus the entropy rate of  
680 the mobility profile would be increased) with little disruption of the utility of the result (because just one  
681 location was replaced/added, and the system can ignore the result of the associated request, by knowing  
682 it is a fake one). However, it would be easy for an attacker to notice the impossibility of the jump, due  
683 to the recent past history, and ignore the location in New York. This happens when we are considering a  
684 mobility profile, since location profiles by their own just account for the number of visits to each place,  
685 leaving unnoticed this kind of impossible large jumps between locations in a short period of time. We  
686 can devise then a semantics and scale-related problem. What data do we want to preserve? For instance,  
687 if only the work/home locations are the ones to be protected, the perturbation methods should focus on  
688 replacing or adding samples of the same city repeatedly, so that their frequency is comparable to the  
689 one of home/work locations. Since the places could be nearby, it would be more difficult for an attacker  
690 to distinguish among the real and fake ones. However, if we want to preserve in which country the  
691 user is, the perturbation mechanism needs to be more sophisticated to make the attacker believe the user  
692 might be at any of two countries by creating equally believable location profiles. Since we need to keep  
693 the scope of this work bounded, we just analyzed the basic cases and made some considerations about  
694 these interesting questions for further research. To the best of our knowledge, there is only another work  
695 aiming at preserving the information contained in the transitions among locations [33]. In this work the  
696 authors face the privacy-utility trade-off by including a function that reflects the distance among the real  
697 location and the fake one that is calculated by their proposed LPPM. Then, the quality (utility) loss is  
698 calculated by averaging the result of applying function to each location that needs to be disclosed for  
699 the LBS (and thus, that can reach the attacker). The distance function reflects which events have more  
700 impact in quality loss, and also the different quality losses of different outputs (fake locations) for each  
701 target location to be protected. However, they do not consider the cost associated to added locations as  
702 an extension of the replacement case. Besides, they do not consider the semantics of the locations and  
703 transitions among them, although they designed their LPPM to face adversaries who can learn how the  
704 LPPM works at each step (thus, potentially being able to notice also if the LPPM is replacing samples  
705 by other ones far away from the current one). Anyhow, this semantics and correlations-related problem  
706 shows to be crucial when assessing the real privacy obtained by the LPPM, and thus needs to be further  
707 investigated in future works.

## 708 **6. Conclusions and future work**

709 In this work we have analyzed privacy-enhancing mechanisms based on information theory concepts,  
710 such as entropy, applied to locations and mobility profiling scenarios. Starting with synthetic and simple  
711 processes, we have shown that the theory applicable to these low alphabet cardinality, memory-less  
712 processes cannot be directly applied to more complex cases, such as mobility profiles of users. Therefore,

713 the remarkable results obtained in the simpler case get degraded until little privacy enhancement is  
714 observed, unless utility is completely lost.

715 The main reasons leading to these results are the increase in the alphabet cardinality (from a few  
716 categories to hundreds of visited places by a user), and the temporal dependencies introduced by the fact  
717 of considering mobility profiles instead of set of independent samples (location profiles), which leads to  
718 the need of using general privacy metrics, such as the one proposed in this work, based on the information  
719 theory concept of entropy rate, in order to consider the temporal dependencies of the mobility profiles.  
720 Moreover, the probability density function underneath in the mobility profile of a user is highly biased  
721 toward certain frequently visited places, which makes it difficult to hide these locations just by replacing  
722 the rest of samples by random locations.

723 As discussed earlier, careful replacement methods should be studied for these special cases. An  
724 interesting future research line might be to investigate how to replace samples taking into account the  
725 current and past locations, in order to provide reasonable replacements, and to exploit the biases toward  
726 the most visited locations to flatten the probability distribution, since these locations and their visitation  
727 profile are the keys to identify the user behind such profiles. Other interesting aspect to explore is the  
728 usefulness of alternative measures of uncertainty, such as the R enyi entropy and the variance, in order to  
729 assess the privacy of mobility profiles.

### 730 **Acknowledgments**

731 This work is partially supported by the Spanish Ministry of Science and Innovation through  
732 CONSEQUENCE (TEC2010-20572-C02-01/02) project. The work of Das was partially supported  
733 by NSF grants IIS-1404673, CNS-1355505, CNS-1404677, and DGE-1433659. Part of the work by  
734 Rodriguez-Carrion was conducted while she was visiting Computer Science Department at Missouri  
735 University of Science and Technology in 2013-2014.

### 736 **Author Contributions**

737 A. Rodriguez-Carrion, C. Campo and C. Garcia-Rubio participated in the conception and  
738 development of the main idea, motivation and discussion, and contributed mainly in the design of the  
739 experiments and manuscript preparation. D. Rebollo-Monedero, J. Forn e and J. Parra-Arnau actively  
740 participated in the conception and development of many of the conceptual, theoretical and experimental  
741 aspects of the paper, but particularly in the information-theoretic formulation and analysis of the problem  
742 investigated. They also made critical revision of the manuscript at all stages of the preparation. S.K. Das  
743 thoroughly revised the paper and provided useful feedback for its improvement. All authors gave final  
744 approval of the version to be submitted.

### 745 **Conflicts of Interest**

746 The authors declare no conflict of interest.

### 747 **References**

- 748 1. Danezis, G. Introduction to privacy technology. Res. talk, Katholieke Univ. Leuven, Comput.  
749 Secur., Ind. Cryptogr. (COSIC), 2007.
- 750 2. Dingedine, R. Free Haven's anonymity bibliography, 2009.
- 751 3. Song, C.; Qu, Z.; Blumm, N.; Barabási, A.L. Limits of Predictability in Human Mobility. *Science*  
752 **2010**, *327*, 1018–1021.
- 753 4. Jaynes, E.T. On the Rationale of Maximum-Entropy Methods. *Proc. IEEE* **1982**, *70*, 939–952.
- 754 5. Jaynes, E.T. Information Theory and Statistical Mechanics II. *Phys. Review Ser. II* **1957**,  
755 *108*, 171–190.
- 756 6. Parra-Arnau, J.; Rebollo-Monedero, D.; Forné, J. Measuring the privacy of user profiles in  
757 personalized information systems. *Future Gen. Comput. Syst.* **2013**. In press.
- 758 7. Wicker, S.B. The Loss of Location Privacy in the Cellular Age. *Commun. ACM* **2012**, *55*, 60–68.
- 759 8. De Mulder, Y.; Danezis, G.; Batina, L.; Preneel, B. Identification via Location-profiling in GSM  
760 Networks. Proceedings of the 7th ACM Workshop on Privacy in the Electronic Society; ACM:  
761 New York, NY, USA, 2008; WPES '08, pp. 23–32.
- 762 9. Freudiger, J.; Shokri, R.; Hubaux, J.P. Evaluating the Privacy Risk of Location-based Services.  
763 Proceedings of the 15th International Conference on Financial Cryptography and Data Security;  
764 Springer-Verlag: Berlin, Heidelberg, 2012; FC'11, pp. 31–46.
- 765 10. Shokri, R.; Theodorakopoulos, G.; Le Boudec, J.Y.; Hubaux, J.P. Quantifying Location Privacy.  
766 Security and Privacy (SP), 2011 IEEE Symposium on, 2011, pp. 247–262.
- 767 11. Samarati, P.; Sweeney, L. Protecting privacy when disclosing information:  $k$ -Anonymity and its  
768 enforcement through generalization and suppression. Tech. rep., SRI Int., 1998.
- 769 12. Samarati, P. Protecting respondents' identities in microdata release. *IEEE Trans. Knowl. Data*  
770 *Eng.* **2001**, *13*, 1010–1027.
- 771 13. Truta, T.M.; Vinay, B. Privacy protection:  $p$ -Sensitive  $k$ -anonymity property. Proc. Int.  
772 Workshop Priv. Data Manage. (PDM); , 2006; p. 94.
- 773 14. Sun, X.; Wang, H.; Li, J.; Truta, T.M. Enhanced  $p$ -sensitive  $k$ -anonymity models for privacy  
774 preserving data publishing. *Trans. Data Priv.* **2008**, *1*, 53–66.
- 775 15. Machanavajjhala, A.; Gehrke, J.; Kiefer, D.; Venkitasubramanian, M.  $l$ -Diversity: Privacy  
776 beyond  $k$ -anonymity. Proc. IEEE Int. Conf. Data Eng. (ICDE); , 2006; p. 24.
- 777 16. Dwork, C. Differential privacy. In *Encyclopedia of Cryptography and Security*; Springer, 2011;  
778 pp. 338–340.
- 779 17. Ho, S.S.; Ruan, S. Differential privacy for location pattern mining. Proceedings of the 4th ACM  
780 SIGSPATIAL International Workshop on Security and Privacy in GIS and LBS. ACM, 2011, pp.  
781 17–24.
- 782 18. Chen, R.; Acs, G.; Castelluccia, C. Differentially private sequential data publication via  
783 variable-length  $n$ -grams. Proceedings of the 2012 ACM conference on Computer and  
784 communications security. ACM, 2012, pp. 638–649.
- 785 19. Bordenabe, N.E.; Chatzikokolakis, K.; Palamidessi, C. Optimal geo-indistinguishable  
786 mechanisms for location privacy. Proceedings of the 2014 ACM SIGSAC Conference on  
787 Computer and Communications Security. ACM, 2014, pp. 251–262.

- 788 20. Chatzikokolakis, K.; Palamidessi, C.; Stronati, M. A predictive differentially-private mechanism  
789 for mobility traces. *Privacy Enhancing Technologies*. Springer, 2014, pp. 21–41.
- 790 21. Shokri, R.; Theodorakopoulos, G.; Papadimitratos, P.; Kazemi, E.; Hubaux, J.P. Hiding in the  
791 Mobile Crowd: Location Privacy through Collaboration. *Dependable and Secure Computing*,  
792 *IEEE Transactions on* **2014**, *11*, 266–279.
- 793 22. Olteanu, A.M.; Huguenin, K.; Shokri, R.; Hubaux, J.P. Quantifying the effect of co-location  
794 information on location privacy. *Privacy Enhancing Technologies*. Springer, 2014, pp. 184–203.
- 795 23. Deng, M. Privacy Preserving Content Protection. PhD thesis, Katholieke Univ. Leuven, Dept.  
796 Elect. Eng. (ESAT), 2010.
- 797 24. Duckham, M.; Mason, K.; Stell, J.; Worboys, M. A formal approach to imperfection in  
798 geographic information. *Comput., Environ., Urban Syst.* **2001**, *25*, 89–103.
- 799 25. Duckham, M.; Kulit, L. A Formal Model of Obfuscation and Negotiation for Location Privacy.  
800 Proc. Int. Conf. Pervas. Comput.; Springer-Verlag: Munich, Germany, 2005; Vol. 3468, *Lecture*  
801 *Notes Comput. Sci. (LNCS)*, pp. 152–170.
- 802 26. Ardagna, C.A.; Cremonini, M.; Damiani, E.; S. De Capitani di Vimercati.; Samarati, P. Location  
803 Privacy Protection Through Obfuscation-Based Techniques. Proc. Annual IFIP Working Conf.  
804 Data Appl. Secur.; Springer-Verlag: Redondo Beach, CA, 2007; Vol. 4602, *Lecture Notes*  
805 *Comput. Sci. (LNCS)*, pp. 47–60.
- 806 27. Yiu, M.L.; Jensen, C.S.; Huang, X.; Lu, H. SpaceTwist: Managing the trade-offs among Location  
807 Privacy, Query Performance, and Query Accuracy in Mobile Services. Proc. IEEE Int. Conf.  
808 Data Eng. (ICDE); , 2008; pp. 366–375.
- 809 28. Kuflik, T.; Shapira, B.; Elovici, Y.; Maschiach, A. Privacy preservation improvement by learning  
810 optimal profile generation rate. User Modeling. Springer-Verlag, 2003, Vol. 2702, *Lecture Notes*  
811 *Comput. Sci. (LNCS)*, pp. 168–177.
- 812 29. Elovici, Y.; Glezer, C.; Shapira, B. Enhancing customer privacy while searching for products and  
813 services on the World Wide Web. *Internet Res.* **2005**, *15*, 378–399.
- 814 30. Shapira, B.; Elovici, Y.; Meshiach, A.; Kuflik, T. PRAW – The model for PRivAte Web. *J. Amer.*  
815 *Soc. Inform. Sci., Technol.* **2005**, *56*, 159–172.
- 816 31. Rebollo-Monedero, D.; Forné, J. Optimal Query Forgery for Private Information Retrieval. *IEEE*  
817 *Trans. Inform. Theory* **2010**, *56*, 4631–4642.
- 818 32. Parra-Arnau, J.; Rebollo-Monedero, D.; Forné, J. Optimal Forgery and Suppression of Ratings  
819 for Privacy Enhancement in Recommendation Systems. *Entropy* **2014**, *16*, 1586–1631.
- 820 33. Theodorakopoulos, G.; Shokri, R.; Troncoso, C.; Hubaux, J.P.; Le Boudec, J.Y. Prolonging the  
821 Hide-and-Seek Game: Optimal Trajectory Privacy for Location-Based Services. Proceedings of  
822 the 13th Workshop on Privacy in the Electronic Society. ACM, 2014, pp. 73–82.
- 823 34. Rebollo-Monedero, D.; Parra-Arnau, J.; Diaz, C.; Forné, J. On the Measurement of Privacy as  
824 an Attacker’s Estimation Error. *Int. J. Inform. Secur.* **2013**, *12*, 129–149.
- 825 35. Cover, T.M.; Thomas, J.A. *Elements of Information Theory*, second ed.; Wiley: New York, 2006.
- 826 36. Serjantov, A.; Danezis, G. Towards an Information Theoretic Metric for Anonymity. Proc.  
827 Workshop Priv. Enhanc. Technol. (PET). Springer-Verlag, 2002, Vol. 2482, pp. 41–53.

- 828 37. Díaz, C.; Seys, S.; Claessens, J.; Preneel, B. Towards measuring anonymity. Proc. Workshop  
829 Priv. Enhanc. Technol. (PET). Springer-Verlag, 2002, Vol. 2482, *Lecture Notes Comput. Sci.*  
830 (*LNCS*), pp. 54–68.
- 831 38. Díaz, C. Anonymity and Privacy in Electronic Services. PhD thesis, Katholieke Univ. Leuven,  
832 2005.
- 833 39. Oganian, A.; Domingo Ferrer, J. A posteriori disclosure risk measure for tabular data based on  
834 conditional entropy. *SORT. 2003, Vol. 27, Núm. 2 [July-December] 2003*.
- 835 40. Voulodimos, A.S.; Patrikakis, C.Z. Quantifying privacy in terms of entropy for context aware  
836 services. *Identity in the Information Society* **2009**, 2, 155–169.
- 837 41. Alfalayleh, M.; Brankovic, L. Quantifying Privacy: A Novel Entropy-Based Measure of  
838 Disclosure Risk. *arXiv preprint arXiv:1409.2112* **2014**.
- 839 42. Rebollo-Monedero, D.; Forné, J.; Domingo-Ferrer, J. From  $t$ -Closeness-Like Privacy to  
840 Postrandomization via Information Theory. *IEEE Trans. Knowl. Data Eng.* **2010**,  
841 22, 1623–1636.
- 842 43. Li, N.; Li, T.; Venkatasubramanian, S.  $t$ -Closeness: Privacy beyond  $k$ -anonymity and  $l$ -diversity.  
843 Proc. IEEE Int. Conf. Data Eng. (ICDE); , 2007; pp. 106–115.
- 844 44. Parra-Arnau, J.; Rebollo-Monedero, D.; Forné, J. A Privacy-Preserving Architecture for the  
845 Semantic Web based on Tag Suppression. Proc. Int. Conf. Trust, Priv., Secur., Digit. Bus.  
846 (TRUSTBUS); , 2010; pp. 58–68.
- 847 45. Parra-Arnau, J.; Rebollo-Monedero, D.; Forné, J.; Muñoz, J.L.; Esparza, O. Optimal tag  
848 suppression for privacy protection in the semantic Web. *Data, Knowl. Eng.* **2012**, 81–82, 46–66.
- 849 46. Gonzalez, M.C.; Hidalgo, C.A.; Barabasi, A.L. Understanding individual human mobility  
850 patterns. *Nature* **2008**, 453, 779–782.
- 851 47. de Montjoye, Y.A.; Hidalgo, C.A.; Verleysen, M.; Blondel, V.D. Unique in the Crowd: The  
852 privacy bounds of human mobility. *Scientific Reports* **2013**, 3.
- 853 48. Hildebrandt, M.; Backhouse, J.; Andronikou, V.; Benoist, E.; Canhoto, A.; Diaz, C.; Gasson, M.;  
854 Geradts, Z.; Meints, M.; Nabeth, T.; Bendegem, J.P.V.; der Hof, S.V.; Vedder, A.; Yannopoulos,  
855 A. Descriptive analysis and inventory of profiling practices – Deliverable 7.2. Technical report,  
856 Future Identity Inform. Soc. (FIDIS), 2005.
- 857 49. Hildebrandt, M.; Gutwirth, S., Eds. *Profiling the European Citizen: Cross-Disciplinary*  
858 *Perspectives*; Springer-Verlag, 2008.
- 859 50. Rodriguez-Carrion, A.; Garcia-Rubio, C.; Campo, C.; Cortés-Martín, A.; Garcia-Lozano, E.;  
860 Noriega-Vivas, P. Study of LZ-Based Location Prediction and Its Application to Transportation  
861 Recommender Systems. *Sensors* **2012**, 12, 7496–7517.
- 862 51. Schurmann, T.; Grassberger, P. Entropy estimation of symbol sequences. *Chaos: An*  
863 *Interdisciplinary Journal of Nonlinear Science* **1996**, 6, 414–427.
- 864 52. Eagle, N.; Pentland, A.; Lazer, D. Inferring Social Network Structure using Mobile Phone Data.  
865 *Proceedings of the National Academy of Sciences (PNAS)* **2009**, 106, 15274–15278.

**Proof.** For each  $t$  (with  $p(t) > 0$ ) and each  $s$ ,

$$p_{S'|T}(s|t) = (1 - \rho) p_{S|T}(s|t) + \rho \frac{1}{k},$$

where  $k$  is the cardinality of the alphabet of  $S$ . Due to the concavity of the entropy and the fact that uniform distributions maximize it, for all  $t$ ,

$$H(S'|t) \geq (1 - \rho) H(S|t) + \rho \log k \geq H(S|t),$$

867 where  $H(S|t)$  denotes the entropy of  $S$  given  $T = t$ , and similarly for  $S'$ . Taking expectations on  $t$ ,  
 868  $H(S'|T) \geq H(S|T)$ . Clearly, equality holds only when  $\rho = 0$ , or else, when  $S$  given  $t$  is uniformly  
 869 distributed, regardless of  $t$ , i.e.,  $p(s|t) = \frac{1}{k} = p(s)$ .  $\square$

## 870 B. Proof of Theorem 2

**Proof.** We prove the statement for the nontrivial case when  $\rho > 0$ . In Lemma 1, take  $S = X_0$ ,  $S' = X'_0$  and  $T = (X_{-1}, \dots, X_{-m})$ , thus

$$H(X'_0|X_{-1}, \dots, X_{-m}) \geq H(X_0|X_{-1}, \dots, X_{-m}),$$

with equality if and only if  $X_0$  is uniform and independent of  $(X_{-1}, \dots, X_{-m})$ . Next, observe that  $X'_0$  and  $(X'_{-1}, \dots, X'_{-m})$  are conditionally independent given  $(X_{-1}, \dots, X_{-m})$ . Apply the conditional-entropy form of the data processing inequality to write

$$H(X'_0|X'_{-1}, \dots, X'_{-m}) \geq H(X'_0|X_{-1}, \dots, X_{-m}),$$

871 with equality if and only if  $X'_0$  and  $(X_{-1}, \dots, X_{-m})$  are conditionally independent given  
 872  $(X'_{-1}, \dots, X'_{-m})$ . Combine both inequalities to immediately conclude the assertions in the theorem  
 873 regarding  $m$  past samples. The claims on the limit of  $m$  for entropy rates follow the same proof, with  
 874  $S = X_0$ ,  $S' = X'_0$  and  $T = (X_{-1}, X_{-2}, \dots)$ .  $\square$

## 875 C. Proof of Theorem 3

**Proof.** In uniform replacement, a sample  $X_n$  will be effectively perturbed when replacement occurs, with probability  $\rho$ , and the replacement sample  $U_n$  does not match the original one. Precisely,

$$\delta = P\{X_n \neq X'_n\} = \rho(1 - P\{X_n = U_n\}).$$

Because  $X_n$  and  $U_n$  are independent and  $U_n$  is uniform,

$$P\{U_n = X_n\} = E_{X_n} P\{U_n = X_n|X_n\} = 1/|\mathcal{X}|.$$

876 If the original process  $X$  is not independent, uniformly distributed, all samples must be replaced to make  
 877 it so, thereby maximizing the entropy rate. Consequently,  $\rho_{\text{crit}} = 1$ , and  $\delta_{\text{crit}}$  can be obtained from the  
 878 relationship between  $\rho$  and  $\delta$  above, simply by setting  $\rho = 1$ .

As for improved replacement, we resort to Theorem 2 in [31] and the concept of critical redundancy, which takes on the value  $1 - \frac{1}{|\mathcal{X}|p_{\max}}$  in the notation of this work. According to this, for any  $\rho \geq 1 - \frac{1}{|\mathcal{X}|p_{\max}}$ , the PMF of the replaced samples  $R_n$  is

$$r(x) = \frac{1}{\rho} \frac{1}{|\mathcal{X}|} + \left(1 - \frac{1}{\rho}\right) p(x).$$

Proceeding as in the first part of this proof,

$$\delta = \rho(1 - \text{P}\{X_n = R_n\}),$$

but now

$$\text{P}\{X_n = R_n\} = \sum_x p(x) r(x),$$

879 from which the expression for  $\delta$  in the second part of the theorem follows.

880 For i.i.d. processes, the problem is mathematically equivalent to that formulated in [31], and  $\rho_{\text{crit}}$   
881 becomes the critical redundancy defined shortly before Theorem 2 in the cited work, in the form  
882 expressed in the statement of the theorem we prove here.

883 The case for processes with memory requires complete replacement to achieve independence of the  
884 samples, not merely uniform distribution, just as in the case of uniform replacement. But for  $\rho = 1$ , the  
885 replacement strategy  $R_n$  becomes uniform, and the analysis for uniform replacement above applies here  
886 as well.  $\square$

887 © May 21, 2015 by the authors; submitted to *Entropy* for possible open access  
888 publication under the terms and conditions of the Creative Commons Attribution license  
889 <http://creativecommons.org/licenses/by/4.0/>.