

# Measuring Re-identification Risk

CJ Carey Google cjcarey@google.com	Travis Dick Google tdick@google.com	Alessandro Epasto Google aepasto@google.com	Adel Javanmard USC*, Google adeljavanmard@google.com
Josh Karlin Google jkarlin@google.com	Shankar Kumar Google shankarkumar@google.com	Andrés Muñoz Medina Google ammedina@google.com	
Vahab Mirrokni Google mirrokni@google.com	Gabriel Henrique Nunes UFMG <sup>†</sup> , Google ghn@nunesgh.com	Sergei Vassilvitskii Google sergeiv@google.com	
	Peilin Zhong Google peilinz@google.com		
	August 1, 2023		

## Abstract

Compact user representations (such as embeddings) form the backbone of personalization services. In this work, we present a new theoretical framework to measure re-identification risk in such user representations. Our framework, based on hypothesis testing, formally bounds the probability that an attacker may be able to obtain the identity of a user from their representation. As an application, we show how our framework is general enough to model important real-world applications such as the Chrome’s Topics API for interest-based advertising. We complement our theoretical bounds by showing provably good attack algorithms for re-identification that we use to estimate the re-identification risk in the Topics API. We believe this work provides a rigorous and interpretable notion of re-identification risk and a framework to measure it that can be used to inform real-world applications.

## 1 Introduction

From curated travel suggestions to localized search results and relevant ads, personalization of online content has become increasingly important, as users expect online systems to be intelligent enough to anticipate their needs. The majority of these systems are now powered by so-called user representations — for instance dense vector embeddings in  $\mathbb{R}^d$  or discrete tokens—that enable classifiers and recommendation systems to generate useful personalized content.

User representations form a compact description of a user profile that distill a user’s interest into a short vector. For instance, a music streaming service may represent a user by a genre of songs the user likes. More sophisticated embeddings may be trained using neural networks and not be as readily explainable.

The compact representation of user profiles may be seen as providing some privacy: instead of describing a user by all raw data (e.g. all songs they have listened to), one can summarize their profile in a few bits.

\*Also affiliated with Data Sciences and Operations, University of Southern California.

<sup>†</sup>Also affiliated with Universidade Federal de Minas Gerais.

Of course, such informal privacy arguments can be incorrect, and do not account for sophisticated attacks or non-obvious data leak vectors. In this work we study the question formally. Specifically, we ask to what extent user representations can be used to link back to and re-identify individuals. We provide a rigorous and interpretable notion of re-identification risk and a framework to measure it.

**Assessing Privacy** Since the seminal work of Dwork et al. [19] introducing differential privacy, significant effort has been devoted into developing differentially private algorithms for numerous problems and settings. Despite these efforts, far from all production systems in use currently employ differential privacy. Moreover, the ever evolving legal framework for privacy (such as GDPR [13] and other regulations) has introduced over time even more privacy definitions, while research has pointed out that these definitions are often not compatible with basic privacy properties like composability [13].

An immediate question is how to assess the privacy of such real-world systems. If the underlying method is randomized, one can try to establish its sensitivity to a single user’s input and derive differential privacy bounds. However, just as in the GDPR example highlighted by [13], the bounds may be vacuous. In this work we propose a complementary approach. We develop a rigorous new framework, centered around re-identification risk, and prove two kinds of results. First, is a series of unconditional lower bounds, showing the risk present in the system. Second, we give a formalization of upper bounds under a “closed-world” assumption, where we can characterize all of the information an attacker has access to. The closed world assumption is natural in some scenarios, and unrealistic in others, in either case our framework can provide a spectrum of the re-identification risk incurred by the users.

We remark that measuring re-identification risk complements the guarantees provided by other notion of privacy, such as (local) differential privacy and k-anonymity, by giving guarantees in a different domain. As part of our analysis, we derive re-id risk guarantees for differentially private algorithms, and k-anonymous datasets. Critically, however, this notion can provide a measure of risk for systems adopting other ad-hoc privacy approaches, an important contribution given the large number of products not designed with differential privacy guarantees.

**Applications** One concrete example we study is the the Topics API [26] proposed by Chrome as part of the Privacy Sandbox initiative [25]. In a nutshell, the Topics API creates a representation of a user corresponding to the top interests of the user in a taxonomy of about 350 interests. When a website calls the Topics API, it returns one of the top interests of the user uniformly at random (and with some small probability a random topic from the taxonomy not necessarily in the user representation). We present a formal study quantifying the extent to which the Topics API can be used to re-identify users across different domains.

We couple this work with an experiment quantifying the risk in releasing samples of music listened by users in a large song dataset. Using our framework, we can quantify re-identification risk in this example, and show that with 4 independent samples per user, the re-id risk is around 1%.

## 1.1 Main contributions

In summary, we make the following main contributions:

- We describe a hypothesis testing [10] framework for evaluating the risk of identifying a random user given access to a representation as well as knowledge of the (potentially randomized) process that generates the representation.
- We extend this scenario and consider a case where an attacker not only observes the representation of a random user but also observes representations of all users and uses this information to match representations to user identities. We refer to this scenario as the matching setting.
- We measure the re-identification risk against an attacker that does not fully know the representation generation distribution (e.g., due to uncertainty on the underlying user data) but has a prior over the set of distributions that generated the representations.

- We specialize this scenario to the Topics API for which we provide an in-depth analysis of an attacker’s re-identification ability under this model.
- We validate our results on two datasets. First a publicly available song dataset, and second a simulation of the Topics API based on proprietary data.

## 2 Related Work

Our work builds upon the rich literature on data anonymization [29], statistical privacy [28], representation learning [8], and web fingerprinting and tracking [31] and more generally on privacy preserving data mining and analysis [3, 23, 45].

Here we discuss only elements of these areas that most relate to the problem of user re-identification. In Section 6 we provide a more in-depth study of how our notion of re-identification relates to other privacy definitions.

**Privacy preserving data releases** This privacy framework traditionally models a data curator that wants to release data without leaking individual user information. In the context of sharing user representations, two of the most popular schemes for doing this are  $k$ -anonymity [42] and differential privacy [17]. In  $k$ -anonymity, the data curator ensures that each user representation is shared by at least another  $k$  users. This is normally achieved by using optimization techniques for minimizing error in redacting tabular data [30, 37, 2, 35], as well as clustering algorithms [21], and tree-based methods [32]. Differential privacy and in particular local differential privacy [48], introduce noise into the user representations to bound that information about a user that is leaked by the representation. There is a vast literature on differential privacy [19, 34, 11, 38, 33, 18] (we refer to [17] for a survey). Significant work in differential privacy has been devoted to private representation learning and machine learning [1].

In Section 6, we show that these two notions are theoretically sufficient (but not necessary) to prove low re-identification risk.

**Browser fingerprinting and web tracking** It is well known that data brokers and ad tech providers can use technologies such as third-party cookies to build detailed browsing history profiles of users [7]. With the impending deprecation of third-party cookies by major browsers, significant attention has been focused by the research community on understanding covert ways of tracking. One such method is browser fingerprinting [31], which relies on characteristics of devices including the screen size and operating system versions to uniquely identify a user across the web. Several studies [20, 24, 5] have analyzed how much information (measured by entropy) trackers can gain from such APIs to create a user identifier. The more entropy an API has, the more likely is that it can be used to re-identifying a user. With some few exceptions [36] most of the research work in this area has focused on understanding the re-identification risk of an API over a fixed instant in time. Our work expands this area of research by modeling the information leakage across time as values returned by an API can change (see section 7). Moreover, we believe our notion of re-identification can have a more straight-forward semantic interpretation compared to the number of bits leaked by an API as it directly measures the re-identification risk.

**Information theoretic notions of privacy** Finally, there is a rich body of literature [47, 46, 16] devoted to studying privacy leakage of information release through the lens of channel capacity. In the context of this paper, a channel is an object that takes as input user information and outputs a user representation. The information theoretic view of privacy measures the ability of an attacker with access to the output of a channel (and with knowledge of the channel internal mechanism) to reconstruct the input to the channel. The quantification of this ability traditionally involves analyzing the mutual information [15] between the input and output of the channel and its generalizations [44]. In Section 6 we show how one can derive bounds on re-identification risk using mutual information. Another related work is that of Cohen et al. [13] that

formalizes privacy as a protection against singling out a user. Unlike our paper, it makes certain assumptions on data generation not present in our work (e.g., that data must be i.i.d. from certain distributions).

Finally, we compare our framework to that of quantitative information flow (QIF) [4], a recent generalization of information theory which provides a flexible and semantically sound framework for analyzing the privacy risks of a channel. As we shall see later on, elements of our framework can be seen as a special case of the QIF framework. By specializing it however, we are able to provide tight characterizations of optimal re-identification attacks.

### 3 Model

We consider a universe of  $n$  users indexed by identities in the set  $\mathcal{I} = [n]$  and an attacker whose goal is to re-identify users based on their representations. Our goal is to quantify the extent to which an attacker can re-identify users as a function of how the user representations are constructed. We study this problem in two settings: First, in the random-user setting, one of the  $n$  users is chosen uniformly at random from  $\mathcal{I}$ , that user’s representation is revealed to the attacker, and the attacker attempts to guess the user’s identity. Next, in the matching setting, the attacker observes representations for all  $n$  users, and their goal is to match each of the  $n$  representations to one of the  $n$  user identities. Notice how this modeling can be seen as a form of hypothesis testing, where the hypothesis corresponds to the possible identities of the user.

We begin by formally defining the processes that generate representations for each user. We then describe the formal attack model, which quantifies how attackers interact with these representations.

**Representations** We assume that user representations belong to a representation space  $\mathcal{O}$  with finite cardinality  $|\mathcal{O}| = m$ . For instance,  $\mathcal{O}$  may be the set of all music genres.<sup>1</sup> Note that we do not make any further structural assumptions on  $\mathcal{O}$ . For example,  $\mathcal{O}$  may be a finite set of topics, songs, images, or points in  $\mathbb{R}^d$ .

An important part of our setup is that a user is not associated with a single representation, but rather with a *distribution* over representations. In other words, the process that assigns a user to a representation may be randomized. As a running example, suppose a user is equally interested in *classical* music and *alt-rock*. One way to represent this behavior is as a fractional assignment  $\{0.5, 0.5\}$  to the two genres. Then, whenever a single genre is required, one can select it from the given distribution. (Looking ahead, this fractional assignment will be key for strengthening re-identification protections.)

Formally, we can encode the representation distributions assigned to all  $n$  users using a row-stochastic *representation matrix*  $\mathbf{P} \in [0, 1]^{n \times m}$  where  $\mathbf{P}[i, o]$  is the probability that user  $i \in \mathcal{I}$  has representation  $o \in \mathcal{O}$  (where we slightly abuse notation and use elements of  $\mathcal{O}$  to index columns of  $\mathbf{P}$ ). We write  $\mathbf{P}[i, \cdot]$  to denote the  $i^{\text{th}}$  row of  $\mathbf{P}$ , which is the user  $i$ ’s distribution over  $\mathcal{O}$ .<sup>2</sup>

Finally, in order to reason about potentially randomized algorithms that assign representations to users, we model the process that constructs the representation matrix  $\mathbf{P}$  (i.e., that assigns representations to users) as a distribution  $\mathcal{D}$  over representation matrices in  $[0, 1]^{n \times m}$ . The distribution  $\mathcal{D}$  models both the user behavior informing the representations and algorithms used to construct the assignment of representation distributions to users.

**Attack Model** We begin with nature sampling a single representation matrix  $\mathbf{P}$  from  $\mathcal{D}$ . At that point, attackers will attempt to re-identify users based on representations sampled from  $\mathbf{P}$ . We detail this further below.

The attacker’s goal is to re-identify users based on their representations (instead of using them for the intended use case—e.g., personalization of content). We categorize attackers by the varying degrees of knowledge about the sampled representation matrix  $\mathbf{P}$ .

<sup>1</sup>We make the finite cardinality assumption for simplicity, as it is sufficient to elucidate the main results of the paper and to model the Topics API. We believe however our theory can be extended to dense embeddings as well.

<sup>2</sup>We stress that the set  $\mathcal{O}$  is arbitrary and this allows representing arbitrary discrete distributions, including over high dimensional spaces. We leave generalizing our framework to continuous distributions as a future work.

In the *full-information* setting, we assume that the attacker observes the sampled representation matrix  $\mathbf{P}$ . This corresponds to a powerful attacker with detailed knowledge of all users and the representation generation process  $\mathcal{D}$ . In the *partial-information* setting, we assume a weaker attacker: one that receives a vector containing one representation sampled from each user’s representation distribution:  $W \in \mathcal{O}^n$  where  $W_i \sim \mathbf{P}[i, \cdot]$ . This corresponds to a situation where an attacker only learns about  $\mathbf{P}$  as a client consuming representations.

In both scenarios, the attacker uses their knowledge of  $\mathbf{P}$  (either the actual  $\mathbf{P}$  itself or the vector of representations  $W$ ) in order to construct a prediction rule  $\varphi$  that they will use to re-identify or match users. In the full-information setting, we require that  $\varphi$  be  $\mathbf{P}$ -measurable, while in the partial-information setting,  $\varphi$  must be  $W$ -measurable.

**Success Metrics** In the random-user setting, the attacker attempts to re-identify a single randomly chosen user based on a sample of their representation.

Formally, the attacker uses their knowledge of  $\mathbf{P}$  to construct a possibly randomized prediction rule  $\varphi_R : \mathcal{O} \rightarrow \mathcal{I}$ . We define the accuracy random variable as follows: let  $I$  be a uniformly random sample from  $[n]$ ,  $O$  be sampled from  $\mathbf{P}[I, \cdot]$ , and define

$$\text{Acc}_R(\varphi_R) = \mathbb{P}(\varphi_R(O) = I \mid \mathbf{P}),$$

which is the probability that the attacker correctly re-identifies the random user conditioned on the representation matrix  $\mathbf{P}$ .

In the matching setting, the attacker receives a set of representations, one for each user, and their goal is to match them to the user identities.

Formally, the attacker uses their knowledge of  $\mathbf{P}$  to construct a matching rule  $\varphi_M : \mathcal{O}^n \rightarrow \mathcal{I}^n$ . We define the matching accuracy random variable as follows: Sample a permutation  $\pi : [n] \rightarrow [n]$  uniformly at random, a vector of independent representations  $O_{1:n} \in \mathcal{O}^n$ , where  $O_i \sim \mathbf{P}[\pi(i), \cdot]$  is a representation sampled for user  $\pi(i)$ , and define

$$\text{Acc}_M(\varphi_M) = \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{\varphi_M^i(O_{1:n}) = \pi(i)\} \mid \mathbf{P} \right],$$

where  $\varphi_M^i(O_{1:n})$  denotes the  $i^{\text{th}}$  component of  $\varphi_M(O_{1:n})$ , which is the attacker’s prediction for the user identity that produced observation  $O_i$ . This is the expected fraction of users the attacker correctly re-identifies, conditioned on the representation matrix  $\mathbf{P}$ .

We do not require that  $\varphi_M(O_{1:n})$  is a permutation of  $\mathcal{I}$ . Note that the random permutation  $\pi$  is used to ensure that the indices in the representation vector  $O_{1:n}$  are not useful for re-identification. In particular, any constant prediction rule  $\varphi_M$  has  $\text{Acc}_M(\varphi_M) = 1/n$ .

## 4 Random-user Accuracy Bounds

In this section we provide information theoretic bounds on the accuracy that any attacker can achieve in the random-user setting. Recall that, in the random user setting, we sample a representation matrix  $\mathbf{P}$  from the distribution  $\mathcal{D}$  and the attacker formulates a (possibly randomized) prediction rule  $\varphi_R : \mathcal{O} \rightarrow \mathcal{I}$  based on their knowledge of  $\mathbf{P}$ . We let  $I$  be a user drawn uniformly at random from  $\mathcal{I}$  and  $O \in \mathcal{O}$  be a representation sampled from  $\mathbf{P}[I, \cdot]$ , and the accuracy random variable,  $\text{Acc}_R(\varphi_R)$ , is the probability that  $\varphi(O) = I$ , conditioned on the representation matrix  $\mathbf{P}$ . Our main result holds with probability one over the draw of the representation matrix  $\mathbf{P}$  from  $\mathcal{D}$ . First, for any prediction rule  $\varphi_R : \mathcal{O} \rightarrow \mathcal{I}$ , we provide an exact expression for  $\text{Acc}_R(\varphi_R)$ . Next, we prove an upper bound on  $\text{Acc}_R(\varphi_R)$  that depends only on  $\mathbf{P}$ . Our upper bound is tight in the full information setting.

Before stating our results, we introduce a matrix representation for any (possibly randomized) prediction rule  $\varphi_R : \mathcal{O} \rightarrow \mathcal{I}$ . Define the matrix  $\mathbf{A} \in [0, 1]^{n \times m}$  to have entries  $\mathbf{A}[i, o] = \mathbb{P}(\varphi_R(o) = i \mid \mathbf{P})$ . That is, conditioned on  $\mathbf{P}$ ,  $\mathbf{A}[i, o]$  is the probability that the attacker’s rule predicts user  $i \in \mathcal{I}$  after receiving representation  $o \in \mathcal{O}$ . With this, we are ready to state our main results.

**Lemma 1.** Let  $\mathbf{P} \sim \mathcal{D}$  be a random representation matrix and  $\mathbf{A}$  be the matrix representation of an attacker's prediction rule  $\varphi_R : \mathcal{O} \rightarrow \mathcal{I}$ . Then, with probability one, we have that

$$\text{Acc}_R(\varphi_R) = \frac{1}{n} \text{tr}(\mathbf{P}\mathbf{A}^\top).$$

*Proof.* Let  $I$  be a user chosen uniformly at random from  $\mathcal{I}$  and  $O \in \mathcal{O}$  be a representation sampled from  $\mathbf{P}[I, \cdot]$ . From the law of total probability, we have that

$$\begin{aligned} \mathbb{P}(\varphi(O) = I | \mathbf{P}) &= \sum_{o \in \mathcal{O}} \sum_{i=1}^n \mathbb{P}(\varphi(O) = I | O = o, I = i, \mathbf{P}) \cdot \mathbb{P}(O = o, I = i | \mathbf{P}) \\ &= \frac{1}{n} \sum_{o \in \mathcal{O}} \sum_{i=1}^n \mathbf{A}[i, o] \cdot \mathbf{P}[i, o], \end{aligned}$$

where the double sum is equal to the Frobenius inner product of  $\mathbf{P}$  and  $\mathbf{A}$ , which can also be written as  $\text{tr}(\mathbf{P}\mathbf{A}^\top)$ .  $\square$

Next, we provide upper bounds on  $\text{Acc}_R(\varphi_R)$  in terms of the representation matrix  $\mathbf{P}$ .

**Corollary 1.** Let  $\mathbf{P} \sim \mathcal{D}$  be a random representation matrix and let  $\varphi_R$  be any prediction rule. Then, with probability one, we have that

$$\text{Acc}_R(\varphi_R) \leq \frac{1}{n} \sum_{o \in \mathcal{O}} \max_{i \in \mathcal{I}} P[i, o] =: \frac{1}{n} \|\mathbf{P}\|_{\infty, 1}.$$

Additionally, there exists an attacker in the full-information setting capable of constructing  $\varphi_R^*$  such that  $\text{Acc}_R(\varphi_R^*) = \frac{1}{n} \|\mathbf{P}\|_{\infty, 1}$ .

*Proof.* Let  $\varphi_R : \mathcal{O} \rightarrow \mathcal{I}$  be any attacker prediction rule and let  $\mathbf{A} \in [0, 1]^{n \times m}$  be its matrix representation. By Lemma 1, the following holds with probability one:

$$\begin{aligned} \text{Acc}_R(\varphi_R) &= \frac{1}{n} \sum_{o \in \mathcal{O}} \sum_{i=1}^n \mathbf{A}[i, o] \mathbf{P}[i, o] \leq \frac{1}{n} \sum_{o \in \mathcal{O}} \sum_{i=1}^n \mathbf{A}[i, o] \cdot \max_{j \in [n]} \mathbf{P}[j, o] \\ &= \frac{1}{n} \sum_{o \in \mathcal{O}} \max_{i \in [n]} \mathbf{P}[i, o] = \frac{1}{n} \|\mathbf{P}\|_{\infty, 1}, \end{aligned}$$

where the second last equality follows from the fact that for any  $o \in \mathcal{O}$ , we have  $\sum_{i=1}^n \mathbf{A}[i, o] = 1$ .

Finally, the inequality in the above derivation holds with equality whenever, for each representation  $o \in \mathcal{O}$  and user  $i \in \mathcal{I}$ ,  $\mathbf{A}[i, o] > 0$  implies that  $\mathbf{P}[i, o] = \max_j \mathbf{P}[j, o]$ . An attacker in the full-information setting can design  $\mathbf{A}$  so that this holds, which implies they are able to achieve this accuracy exactly. In other words, if the attacker designs  $\varphi_R^*$  so that  $\varphi_R^*(o)$  is always some user  $i$  with the maximum probability of generating representation  $o$ , then  $\text{Acc}_R(\varphi_R) = \frac{1}{n} \|\mathbf{P}\|_{\infty, 1}$ .  $\square$

Next, we prove a bound on the accuracy of an attacker with partial-information. In particular, we bound their expected accuracy conditioned on the partial information contained in  $W \in \mathcal{O}^n$ . This is the probability that the attacker will correctly predict the identity of a random user based on a sample representation, conditioned on vector  $W$ .

**Lemma 2.** Let  $\mathbf{P} \sim \mathcal{D}$  be a random representation matrix,  $W \in \mathcal{O}^n$  be a vector of independent representations where  $W_i \sim \mathbf{P}[i, \cdot]$ , and let  $\varphi_R : \mathcal{O} \rightarrow \mathcal{I}$  be a prediction rule that is  $W$ -measurable. Then we have that

$$\mathbb{E}[\text{Acc}_R(\varphi_R) | W] \leq \frac{1}{n} \|\mathbb{E}[\mathbf{P} | W]\|_{\infty, 1}.$$

Moreover there exists a prediction rule  $\varphi_R^*$  for which the upper bound is achieved.

*Proof.* Let  $\mathbf{A} \in [0, 1]^{n \times m}$  be the matrix-representation of the attacker's  $W$ -measurable prediction rule  $\varphi_R$ . With probability one, we have that  $\text{Acc}_R(\varphi_R) = \frac{1}{n} \text{tr}(\mathbf{P}\mathbf{A}^\top)$ . From this, it follows that

$$\mathbb{E}[\text{Acc}_R(\varphi_R) | W] = \mathbb{E} \left[ \frac{1}{n} \text{tr}(\mathbf{P}\mathbf{A}^\top) \middle| W \right] = \frac{1}{n} \text{tr}(\mathbb{E}[\mathbf{P} | W]\mathbf{A}^\top),$$

where the final equality follows from the fact that  $\mathbf{A}$  is  $W$ -measurable. Next, since  $\mathbb{E}[\mathbf{P} | W]$  is a row-stochastic matrix, the same argument as in Corollary 1, it follows that  $\frac{1}{n} \text{tr}(\mathbb{E}[\mathbf{P} | W]\mathbf{A}^\top) \leq \|\mathbb{E}[\mathbf{P} | W]\|_{\infty, 1}$ , as required. To find the prediction rule  $\varphi_R^*$  we use the same argument as in Corollary 1 and define a matrix  $\mathbf{A}$  such that  $\mathbf{A}[i, o] > 0$  implies  $i \in \text{argmax}_j \mathbb{E}[\mathbf{P} | W]$ .  $\square$

## 5 Matching model accuracy bounds

In this section, we provide an information theoretic upper bound on the accuracy that any attacker can achieve in the full-information setting. Recall that in the random-user setting, we sample a representation matrix  $\mathbf{P}$  from the distribution  $\mathcal{D}$ , the attacker formulates a (possibly randomized) prediction rule  $\varphi_M : \mathcal{O}^n \rightarrow \mathcal{I}^n$  based on their knowledge of  $\mathbf{P}$ . Then we let  $\pi : [n] \rightarrow [n]$  be a permutation of  $[n]$  chosen uniformly at random,  $O_{1:n} \in \mathcal{O}^n$  be a vector of independent observations with  $O_i \sim \mathbf{P}[\pi(i), \cdot]$ , and define the accuracy random variable by  $\text{Acc}_M(\varphi_M) = \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{\varphi_M^i(O_{1:n}) = \pi(i)\} \middle| \mathbf{P} \right]$ . Our accuracy bound for the matching setting is given below:

**Lemma 3.** *Let  $\mathbf{P} \sim \mathcal{D}$  be a random representation matrix and let  $\varphi_M : \mathcal{O}^n \rightarrow \mathcal{I}^n$  be the matching rule constructed by an attacker. Then with probability one over  $\mathbf{P} \sim \mathcal{D}$ , we have*

$$\text{Acc}_M(\varphi_M) \leq \frac{m}{n} - \frac{1}{n} \sum_{o \in \mathcal{O}} \prod_{i \in \mathcal{I}} (1 - \mathbf{P}[i, o]).$$

*Proof.* Let  $\mathbf{P}$  be sampled from  $\mathcal{D}$  and  $\varphi_M : \mathcal{O}^n \rightarrow \mathcal{I}^n$  be the matching rule constructed by the attacker based on their knowledge of  $\mathbf{P}$ . Next, let  $O_{1:n}^*$  be independent representations with  $O_i^* \sim \mathbf{P}[i, \cdot]$ , let  $\pi : [n] \rightarrow [n]$  be a permutation of  $[n]$  chosen uniformly at random, and define  $O_{1:n}$  by  $O_i = O_{\pi(i)}^*$ . Next, for each observation  $o \in \mathcal{O}$ , let  $S_o = \{i : O_i = o\}$  denote the set of indices  $i$  for which  $O_i = o$ . Since  $\pi$  is uniformly random and  $O_1, \dots, O_n$  are independent with  $O_i \sim \mathbf{P}[\pi(i), \cdot]$ , we have that

$$\begin{aligned} \text{Acc}_M(\varphi_M) &= \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{\varphi_M^i(O_{1:n}) = \pi(i)\} \middle| \mathbf{P} \right] \\ &= \mathbb{E} \left[ \frac{1}{n} \sum_{o \in \mathcal{O}} \sum_{i \in S_o} \mathbb{P}(\varphi_M^i(O_{1:n}) = \pi(i) | O_{1:n}, O_{1:n}^*, \mathbf{P}) \middle| \mathbf{P} \right], \end{aligned}$$

where the second equality follows from breaking the sum over  $i$  into a sum over  $o \in \mathcal{O}$  and  $i \in S_o$ , and adding an inner expectation conditioned on  $O_{1:n}$ ,  $O_{1:n}^*$ , and  $\mathbf{P}$ . The key idea is that, conditioned on  $O_{1:n}$  and  $O_{1:n}^*$ , the permutation  $\pi$  is still random, but  $\varphi_M^i(O_{1:n})$  is fixed, which implies that  $\varphi_M^i(O_{1:n})$  cannot be correct with too large of a probability. For any  $i \in S_o$ , we have that

$$\mathbb{P}(\pi(i) = j | O_{1:n}, O_{1:n}^*, \mathbf{P}) = \frac{\mathbb{I}\{O_j^* = o\}}{|S_o|}.$$

With this, we have

$$\begin{aligned} \text{Acc}_M(\varphi_M) &= \mathbb{E} \left[ \frac{1}{n} \sum_{o \in \mathcal{O}: S_o \neq \emptyset} \sum_{i \in S_o} \frac{\mathbb{I}\{O_{\varphi_M^i(O_{1:n})}^* = o\}}{|S_o|} \middle| \mathbf{P} \right] \\ &\leq \mathbb{E} \left[ \frac{1}{n} \sum_{o \in \mathcal{O}: |S_o| \neq 0} 1 \middle| \mathbf{P} \right] = \mathbb{E} \left[ \frac{1}{n} \sum_o \mathbb{I}\{S_o \neq \emptyset\} \middle| \mathbf{P} \right]. \end{aligned}$$

This final expression is the number of unique representations that the attacker observed, divided by the number of users  $n$ . Intuitively, this bound follows from the fact that for all the users that generated the same observation, the expected number of correct guesses of the attacker is at most one.

To finish the proof, we compute the expected number of distinct representations that the attacker will observe. We have

$$\mathbb{E} \left[ \sum_o \mathbb{I}\{S_o \neq \emptyset\} \mid \mathbf{P} \right] = \sum_o \mathbb{P}(S_o \neq \emptyset \mid \mathbf{P}) = \sum_o \left( 1 - \prod_{i=1}^n (1 - \mathbf{P}[i, o]) \right).$$

which implies the statement of the Lemma.  $\square$

## 5.1 Connections between the random-user and matching models

In this section we study connections between the random-user and matching settings. In particular, we show that the matching setting is no harder for the attacker than the random-user setting: we prove that any attacker in the random-user setting can be modified to achieve the same accuracy in the matching setting. Next, we show that there exist representation probability matrices  $\mathbf{P}$  such that an optimal attacker in the matching setting can do strictly better than the optimal attacker in the random-user setting.

**Lemma 4.** *Let  $\mathbf{P} \sim \mathcal{D}$  and  $\varphi_R : \mathcal{O} \rightarrow \mathcal{I}$  be any (possibly randomized) attacker prediction rule for the random-user setting. Define  $\varphi_M : \mathcal{O}^n \rightarrow \mathcal{I}^n$  by  $\varphi_M(O_{1:n}) = (\varphi_R(O_1), \dots, \varphi_R(O_n))$ . Then with probability one over  $\mathbf{P}$ , we have  $\text{Acc}_M(\varphi_M) = \text{Acc}_R(\varphi_R)$ .*

*Proof.* Intuitively, the matching rule  $\varphi_M$  applies the random-user rule independently for each representation vector in  $O_{1:n}$ , and the expected fraction of entries it will predict correctly is equal to the expected accuracy of  $\varphi_R$  in the random-user setting. Formally, let  $\pi : [n] \rightarrow [n]$  be the random permutation used in the matching setting. Then we have

$$\begin{aligned} \text{Acc}_M(\varphi_M) &= \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{\varphi_M^i(O_{1:n}) = \pi(i)\} \mid \mathbf{P} \right] \\ &= \sum_{i=1}^n \frac{1}{n} \mathbb{P}(\varphi_M^i(O_{1:n}) = \pi(i) \mid \mathbf{P}) \\ &= \sum_{i=1}^n \frac{1}{n} \mathbb{P}(\varphi_R(O_i) = \pi(i) \mid \mathbf{P}) \\ &= \text{Acc}_R(\varphi_R), \end{aligned}$$

where the final equality follows from the fact that the pair  $(I, O)$  with  $I = \pi(i)$  and  $O = O_i$  is distributed identically to  $I \sim \text{Uniform}(n)$  and  $O \sim \mathbf{P}[I, :]$ , since when  $\pi$  is a random permutation,  $\pi(i)$  is randomly chosen uniformly at random from  $[n]$ .  $\square$

Next, we construct a distribution  $\mathcal{D}$  over representation matrices  $\mathbf{P}$  such that an attacker in the matching setting can have a constant factor higher accuracy than the best attacker in the random-user setting.

**Lemma 5.** *For any even number of users  $n$ , there exists a representation space  $\mathcal{O}$  of size  $m = \frac{3n}{2}$  and a distribution  $\mathcal{D}$  over representation matrices  $\mathbf{P} \in \mathbb{R}^{n \times m}$  such that: with probability one, every  $\varphi_R : \mathcal{O} \rightarrow \mathcal{I}$  has  $\text{Acc}_R(\varphi_R) \leq \frac{3}{4}$  in the random-user setting, and there exists a rule  $\varphi_M : \mathcal{O}^n \rightarrow \mathcal{I}^n$  such that  $\text{Acc}_M(\varphi_M) = \frac{7}{8}$  in the matching setting.*

*Proof.* For simplicity, we construct  $\mathcal{D}$  as a distribution supported on a single representation matrix  $\mathbf{P}$ . First, consider the case where we have only  $n = 2$  users, the representation space is  $\mathcal{O} = \{u_1, u_2, a\}$ , and the

representation probability matrix is defined by

$$\mathbf{P} = \begin{array}{c} \text{User 1} \\ \text{User 2} \end{array} \begin{array}{ccc} u_1 & u_2 & a \\ \left( \begin{array}{ccc} 1/2 & 0 & 1/2 \\ 0 & 1/2 & 1/2 \end{array} \right). \end{array}$$

We have  $\frac{1}{n}\|\mathbf{P}\|_{\infty,1} = 3/4$ , and it follows that no prediction rule  $\varphi_R$  can achieve accuracy higher than  $3/4$  in the random-user setting. However, in the matching setting, if the attacker observes at least one of  $\{u_1, u_2\}$ , this is sufficient for perfectly matching the users, since  $u_i$  is only ever generated by user  $i$ . When both users generate the ambiguous representation  $a$ , the attacker still has a  $1/2$  chance to correctly identify the users (e.g., by predicting a random permutation). The probability that both users generate representation  $a$  is  $1/4$ , and the probability that at least one of  $u_1$  or  $u_2$  is generated is  $3/4$ . It follows an attacker in the matching setting can achieve:

$$\begin{aligned} \text{Acc}_M(\varphi_M) &= 1 \cdot \mathbb{P}(u_1 \text{ or } u_2 \text{ observed}) + \frac{1}{2} \cdot \mathbb{P}(\text{only } a \text{ observed}) \\ &= 1 \cdot \frac{3}{4} + \frac{1}{2} \cdot \frac{1}{4} = \frac{7}{8}. \end{aligned}$$

To extend this example to any even number of users, we create  $n/2$  copies of the 2-user problem as follows: Let  $\mathcal{O} = \{u_1, \dots, u_n\} \cup \{a_1, \dots, a_{n/2}\}$  and define the representation probability matrix by

$$\mathbf{P}[i, o] = \begin{cases} 1/2 & \text{if } o = u_i \text{ or } o = a_{\lceil i/2 \rceil}. \\ 0 & \text{otherwise.} \end{cases}$$

Then  $\mathbf{P}$  has  $\frac{3n}{2}$  columns and the maximum value in each column is  $1/2$ . It follows that  $\frac{1}{n}\|\mathbf{P}\|_{\infty,1} = \frac{1}{n} \cdot \frac{3n}{2} \cdot \frac{1}{2} = \frac{3}{4}$ . On the other hand, for any even index  $i$ , we know that there are exactly two entries in  $O_{1:n}$  in the set  $\{u_{i-1}, u_i, a_{i/2}\}$ , and that these entries must correspond to users  $i-1$  and  $i$  (but we do not know the order). When the attacker attempts to identify users  $i-1$  and  $i$ , they are faced exactly with the two-user problem described above, and their expected accuracy for users  $i-1$  and  $i$  is  $7/8$ . Averaging over the  $n/2$  pairs of users, their overall accuracy is also  $7/8$ .  $\square$

## 6 Relation to other privacy notions

In this section we give a detailed discussion of the re-identification risk introduced in Section 3 in relation to two prior notions of algorithmic privacy: local differential privacy and  $k$ -anonymity. We show that (for appropriate parameters) both of these privacy notions are *sufficient* to imply low re-identification risk, but neither condition is *necessary* to obtain low re-identification risk in our framework. (We refer however to the discussion in Section 9 on why they may still be needed for other privacy risks.)

This shows that the re-identification risk outlined in Section 3 is not entirely captured by either of these concepts. We conclude this section by discussing the connection between our work and the field of quantitative information flow (QIF) [4].

### 6.1 Local differential privacy

Local differential privacy (LDP) [14] is a strong privacy notion applicable to publishing user representations constructed from private information. Intuitively, it should be hard to derive the identity of a user from the output of a differentially private mechanism. In this section we prove this implication, while, at the same time, showing that local differential privacy is *not necessary* for low re-identification risk. This result highlights the ability of our framework to characterize directly and sharply re-identification risks.

**Definition 1** (Local differential privacy). Let  $\mathcal{X}$  be an arbitrary space encoding user information. A randomized algorithm  $\mathcal{A} : \mathcal{X} \rightarrow \mathcal{O}$  for mapping user data to a representation satisfies  $(\epsilon, \delta)$ -LDP if the following holds: for all  $x, x' \in \mathcal{X}$  and any set of representations  $E \subset \mathcal{O}$ , we have that

$$\mathbb{P}(\mathcal{A}(x) \in E) \leq e^\epsilon \cdot \mathbb{P}(\mathcal{A}(x') \in E) + \delta.$$

LDP representations as described above can be modeled under our framework as follows: Let  $\mathcal{D}$  be a distribution over representation matrices that samples  $\mathbf{P}$  in two steps: First, the  $n$  users generate their data  $x_1, \dots, x_n \in \mathcal{X}$ . Second, we define  $\mathbf{P} \in [0, 1]^{n \times m}$  to have entries  $\mathbf{P}[i, o] = \mathbb{P}(\mathcal{A}(x_i) = o)$ , where the probability is only over the randomness of the mechanism  $\mathcal{A}$ . Sampling  $\mathbf{P}$  from  $\mathcal{D}$  corresponds to the process generating the users data, while the matrix  $\mathbf{P}$  encodes the mechanism  $\mathcal{A}$ 's output distribution for each user. Since  $\mathcal{A}$  is  $(\epsilon, \delta)$ -LDP, with probability one over the draw of  $\mathbf{P}$ , we have that for any users  $i$  and  $j$ , and any representation subset  $E \subset \mathcal{O}$ , the following holds:

$$\sum_{o \in E} \mathbf{P}[i, o] \leq \delta + e^\epsilon \cdot \sum_{o \in E} \mathbf{P}[j, o]. \quad (1)$$

More generally, we say that any distribution  $\mathcal{D}$  over representation matrices  $\mathbf{P}$  that satisfy (1) with probability one is  $(\epsilon, \delta)$ -LDP.

The following result shows that  $(\epsilon, \delta)$ -LDP implies low re-identification accuracy in the random-user setting.

**Lemma 6.** *Let  $\mathcal{D}$  be any distribution that satisfies  $(\epsilon, \delta)$ -LDP, let  $\mathbf{P} \sim \mathcal{D}$ , and  $\varphi_R$  be an attacker's prediction rule in the random user setting. Then*

$$\text{Acc}_R(\varphi_R) \leq \frac{e^\epsilon + \min(n, m) \cdot \delta}{n}.$$

*Proof.* Let  $\mathbf{P} \sim \mathcal{D}$  be a sampled representation matrix and partition  $\mathcal{O}$  into sets  $\mathcal{O}_1, \dots, \mathcal{O}_n$ , where  $\mathcal{O}_i$  contains all representations that user  $i$  generates with higher probability than any other user with ties broken in favor of the user with lower index. That is,

$$\mathcal{O}_i = \{o \in \mathcal{O} : \text{for all } j \neq i, \mathbf{P}[i, o] \geq \mathbf{P}[j, o] \text{ and if } \mathbf{P}[i, o] = \mathbf{P}[j, o] \text{ then } i < j\}.$$

Then we have that

$$\|\mathbf{P}\|_{\infty, 1} = \sum_{o \in \mathcal{O}} \max_{i \in \mathcal{I}} \mathbf{P}[i, o] = \sum_{i \in \mathcal{I}} \sum_{o \in \mathcal{O}_i} \mathbf{P}[i, o]. \quad (2)$$

Now suppose that  $\mathcal{D}$  is  $(\epsilon, \delta)$ -LDP. Then we have that:

$$\|\mathbf{P}\|_{\infty, 1} = \sum_{o \in \mathcal{O}} \max_{i \in \mathcal{I}} \mathbf{P}[i, o] \leq \sum_{o \in \mathcal{O}} (e^\epsilon \mathbf{P}[1, o] + \delta) = e^\epsilon + m\delta.$$

At the same time, from (2) we have that

$$\|\mathbf{P}\|_{\infty, 1} = \sum_{i \in \mathcal{I}} \sum_{o \in \mathcal{O}_i} \mathbf{P}[i, o] \leq \sum_{i \in \mathcal{I}} \left( \delta + e^\epsilon \sum_{o \in \mathcal{O}_i} \mathbf{P}[1, o] \right) = e^\epsilon + n\delta.$$

The above arguments show that  $\|\mathbf{P}\|_{\infty, 1} \leq e^\epsilon + \min(n, m) \cdot \delta$ . From Corollary 1, it follows that  $\text{Acc}_R(\varphi_R) \leq \frac{e^\epsilon + \min(n, m)\delta}{n}$ .  $\square$

## 6.2 k-anonymity

A process that releases anonymized data about a collection of users is said to be  $k$ -anonymous if the information released for each user cannot be distinguished from at least  $k - 1$  other users who also appear in the release. We say that a distribution  $\mathcal{D}$  over representation matrices is  $k$ -anonymous if, with probability one over  $\mathbf{P} \sim \mathcal{D}$ , every row of  $\mathbf{P}$  is a one-hot vector and appears at least  $k$  times. That is, each user  $i$  is assigned a representation  $o_i \in \mathcal{O}$  that is shared with at least  $k - 1$  other users, and their row of  $\mathbf{P}$  is given by  $\mathbf{P}[i, o] = \mathbb{I}\{o = o_i\}$ .

The following result shows that  $k$ -anonymity is sufficient to limit an attacker’s accuracy to  $1/k$  in the random-user setting.

**Lemma 7.** *Let  $\mathcal{D}$  be any distribution that satisfies  $k$ -anonymity, let  $\mathbf{P} \sim \mathcal{D}$ , and  $\varphi_R$  be an attacker’s prediction rule in the random user setting. Then*

$$\text{Acc}_R(\varphi_R) \leq \frac{1}{k}.$$

*Proof.* Let  $\mathbf{P} \sim \mathcal{D}$  and  $o_1, \dots, o_n \in \mathcal{O}$  be the corresponding representations assigned to each user. For each representation  $o \in \mathcal{O}$ , let  $\mathcal{I}_o = \{i \in \mathcal{I} \mid o_i = o\}$  denote the set of users assigned that representation. From the  $k$ -anonymity condition, we are guaranteed that either  $|\mathcal{I}_o| = 0$  or  $|\mathcal{I}_o| \geq k$ . It follows that there are at most  $n/k$  observations  $o$  for which  $|\mathcal{I}_o| > 0$ . This implies that

$$\|\mathbf{P}\|_{\infty,1} = \sum_{o \in \mathcal{O}} \max_{i \in \mathcal{I}} \mathbf{P}[i, o] = \sum_{o \in \mathcal{O}} \mathbb{I}\{|\mathcal{I}_o| > 0\} \leq \frac{n}{k}.$$

By Corollary 1, it follows that  $\text{Acc}(\varphi_R) \leq \frac{1}{k}$ , as required. □

## 6.3 LDP and k-anonymity are not necessary conditions for low re-identification risk

In the previous two subsections we showed that, for appropriate parameter settings,  $(\epsilon, \delta)$ -LDP and  $k$ -anonymity both imply that an attacker in the random-user setting has low accuracy. In this section, we show that neither condition is necessary.

**Lemma 8.** *There exist distributions  $\mathcal{D}$  such that with probability one, every attacker has  $\text{Acc}_R(\varphi_R) \leq \frac{2}{n}$  and  $\mathcal{D}$  is not  $(\epsilon, \delta)$ -LDP unless  $\delta = 1$  and not  $k$ -anonymous for any  $k$ .*

*Proof.* Let  $\mathcal{O} = \{1, 2\}$  and let  $\mathbf{P}$  have entries given as follows: for each user  $i \in [n]$ , define

$$\mathbf{P}[i, 1] = 1 - \frac{(i-1)}{n-1} \quad \text{and} \quad \mathbf{P}[i, 2] = \frac{(i-1)}{n-1}.$$

Not all of the rows of  $\mathbf{P}$  are one-hot, so it does not satisfy the  $k$ -anonymity requirement. Next, we have that  $\mathbf{P}[1, 2] = 0$  while  $\mathbf{P}[n, 2] = 1$ , which implies that  $\mathbf{P}$  only satisfies the  $(\epsilon, \delta)$ -LDP constraint when  $\delta = 1$ . Finally, we have that  $\|\mathbf{P}\|_{\infty,1} = 2$  and by Corollary 1 it follows that  $\text{Acc}_R(\varphi_R) \leq \frac{2}{n}$ . □

## 6.4 Mutual Information

Another view of re-identification risk can be obtained from the field of information theory. Given a joint pair of random variables  $(X, Y)$ , we are interested in measuring how much *information* does  $Y$  encode about  $X$ . For our random user model, this can be translated to measuring how much information the representation  $O$  provides about the identity random variable  $I$ . This concept is formalized by the conditional mutual information [15]  $\text{MI}(I; O|\mathbf{P})$  for the full information scenario and by  $\text{MI}(I; O|W)$  for the partial-information scenario. This metric was in fact used in prior work [22] to quantify the re-identification risk of the Topics API. One can use the celebrated Fano’s inequality [15] to show that.

**Lemma 9.** *Under the random-user model we have*

$$\text{Acc}_R(\varphi_R) \leq \frac{1 + MI(I; O | \mathbf{P})}{\log(n)}, \quad (3)$$

It is worth noticing that the dependency on the number of users  $n$  here is logarithmic as opposed to that of Lemma 1 where the dependency is linear. This is an exponential improvement and demonstrates that our framework can better capture re-identification risks.

## 6.5 Quantitative information flow

Quantitative information flow [40, 4] (QIF) is a different framework for analyzing the privacy vulnerability of a system. QIF is specified by a space of secrets  $\mathcal{S}$ , an output space  $\mathcal{O}$ , a (possibly randomized) channel  $C: \mathcal{S} \rightarrow \mathcal{O}$  assumed to be known to an adversary, and a gain function  $g: \mathcal{W} \times \mathcal{S} \rightarrow \mathbb{R}$ , where  $\mathcal{W}$  is an adversary’s space of strategies that may coincide with  $\mathcal{S}$  depending on the scenario.

QIF assumes there is a secret  $s$  sampled from a known distribution  $\pi$  and that the adversary observes  $o = C(s)$ . The goal of the adversary is, given  $o$ , to learn about  $s$ . The gain function may then capture the reward of an adversary predicting secret  $s'$ , here a reward function  $r(s', s)$ . Note that an adversary with access to the channel, given an output  $o$ , can always predict the secret  $s'$  that maximizes their posterior reward

$$R(o) = \max_{s'} \sum_{s \in \mathcal{S}} r(s', s) P(s|o).$$

The privacy vulnerability of a channel may be seen in QIF as the expected posterior reward  $\mathbb{E}_\pi[R(o)]$ . For our full-information setting, the known representation matrix  $\mathbf{P}$  corresponds to the channel, the identity space  $\mathcal{I}$  is the secret space and the reward function  $r(s', s) = 1$  if  $s' = s$  and it is 0 otherwise. That is, the adversary is only rewarded if they predict the correct user in one try. It is known [4] that for this scenario the expected posterior reward corresponds to the so-called Bayes vulnerability and it satisfies:

$$\mathbb{E}[R(o)] = \frac{1}{n} \|\mathbf{P}\|_{\infty, 1}$$

That is, our full-information setting is an alternative formulation of QIF as a hypothesis testing framework. To the best of our knowledge, we are not aware of a partial information QIF formulation that fully matches the random-user or matching scenarios although we are actively exploring ways to use advanced concepts in QIF such as the internal fixed-probability choice model [4] to establish a similar connection.

## 7 Case Study: The Topics API

As mentioned in the introduction, we will use our framework to provide an analysis of the re-identification risk in context of the Topics API [26] of the Privacy Sandbox [25]. Here, we introduce the Topics API using the framework of section 3.

The Privacy Sandbox [25] is a series of proposals to enable online advertising while limiting cross-site tracking on the web. We will focus on Interest Based Advertising (IBA) use case of the Privacy Sandbox. IBA is a sector of online advertising in which ad-tech providers build models of the users’ interests in an effort to show them relevant ads. For instance, people interested in a car may be served car ads even on unrelated pages.

Historically, IBA has been enabled through third-party cookies. These serve as a cross-site user identifier, allowing ad techs to keep track of the sites a user has visited and build an interest profile based on their browsing history. This cross-site tracking is in direct contrast with the goals of the Privacy Sandbox, which has led Chrome to announce the Topics API to support IBA without relying on cross-site tracking.

The Topics API works as follows (we refer to the specifications in [26]): every week the *browser* builds an interest profile of the user, in the form of selecting top five topics from a fixed topics taxonomy,  $\mathcal{T}$ . Importantly, this profile is kept on the browsers and is not shared with others.

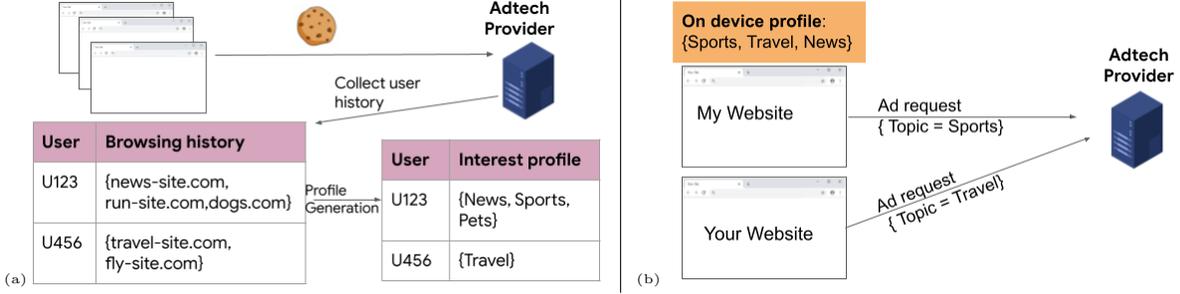


Figure 1: Comparison of IBA under (a) third party cookies and (b) the Topics API. In the former, the ad tech provider gets to build the browsing history of a user. In the latter, the ad tech provider only observes a single topic for this user.

Whenever a website wants to show an ad, the browser shares a topic selected uniformly at random from one of the top 5 topics in the profile of the previous week with the ad tech provider (additionally with some probability  $p$ , it may simply return a uniformly at random topic from  $\mathcal{T}$ ). Crucially, for every user, the topic sampled for a website is fixed for an entire week, and the samples on two different websites are independent.<sup>3</sup> See Figure 1 for an example of the Topics API. A detailed specification of the API can also be found in Algorithm 1.

---

**Algorithm 1** Topics API.

**Input:** Topics  $\mathcal{T}$ , probability to return a random topic  $p$ .

---

**On device:** Select set  $S$  of top 5 most popular topics for this client.

**On call GetTopic()** from website  $w$  on week  $s$  and user  $u$ :

Seed the random number generator with  $w, s, u$ .

Flip coin with heads probability  $p$ .

**if** Heads **then**

**return** Element of  $\mathcal{T}$  chosen u.a.r.

**else**

**return** Element of  $S$  chosen u.a.r.

**end if**

---

**Re-identification risks of the Topics API** Compared to third-party cookies, the Topics API has a significantly lower risk of re-identifiability (the former guarantees re-identifiability by nature of being a persistent cross-site identifier). The goal of this section is to formally measure this risk in the case of misuse of the API.

Consider a user that visits two sites every week. Over time, even with the randomization, the sequence of topics observed on website 1 will be similar to the sequence of topics observed on website 2 (for instance, we expect an exact match in  $(1 - p)/5$  fraction of the weeks). Thus two sites could try to collude to use the Topics API to link the identity of a user across them.<sup>4</sup>

Figure 2 shows how this attack may happen. We now formalize the re-identification risks of the Topics API from the perspective of website 2 as an attacker colluding with website 1.

Let  $N = |\mathcal{T}|$  and  $\mathcal{O} = \mathcal{T}^r$ . That is if  $o \in \mathcal{O}$  then  $o = (o^1, \dots, o^r)$  where  $o^s$  corresponds to the topic returned by the API on week  $s$ . We begin by modeling a single representation of the Topics API. Based on

<sup>3</sup>The API actually returns, on top of the current sampled topic, a cached result for the output of the previous 2 weeks for the caller. We omit this detail from the modeling as observing  $r$  consecutive weeks of Topics, simply corresponds to performing calls for  $r + 2$  weeks in our model.

<sup>4</sup>We restrict our analysis to two sites only for simplicity of exposition, as two sites are sufficient to elucidate the re-identification risk. We refer to section 9 for the limitations of our work.

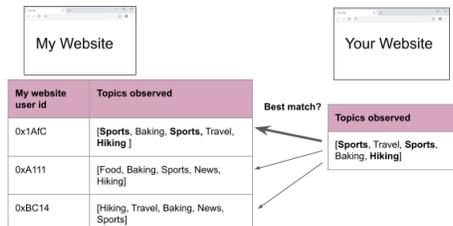


Figure 2: Example of two sites using the sequence of topics returned by the API to potentially re-identify a user across two sites. In the example the user on website 2 shares 3 topics with a user on website 1. This has the most matches across all three users.

Algorithm 1, for every user  $i \in \mathcal{I}$ , the topic returned by the API at round  $s$  depends on the set of top topics  $S_i^s$  associated with the user. Having fixed this set, the topic selection of the API at round  $s$  can be modeled by a matrix  $\mathbf{P}_s$  given by:

$$\mathbf{P}_s[i, o^s] = \begin{cases} q_{\text{in}} := (1 - p)/5 + p/N & o^s \in S_i^s \\ q_{\text{out}} := p/N & o^s \notin S_i^s \end{cases} \quad (4)$$

It is easy to see that the representation matrix for  $o \in \mathcal{O}$  is given by

$$\mathbf{P}[i, o] = \prod_{s=1}^r \mathbf{P}[i, o^s]$$

Thus far we have assumed the top set of topics  $S_i^s$  to be fixed. However, notice that these are not observable by the websites directly, since they can only see the samples from the top sets. Moreover, the sets  $S_i^s$  are determined by the behaviors of the users which we encode as a probabilistic process in the distribution  $\mathcal{D}$ .

Formally, the matrix  $\mathbf{P}$  is a random variable sampled from a (latent) distribution  $\mathcal{D}$  which encodes the way that the top set of topics  $S_i^s$  is generated. Clients of the Topics API (i.e., websites) learn about the sampled matrix  $\mathbf{P}$  through the observation of samples of the representations in their site.

With this modeling in mind, we see that this re-identification risk can be defined in terms of the (partial) information random-user model:

1. A matrix  $\mathbf{P}$  is sampled from  $\mathcal{D}$ .
2. A sample of representations  $W = (W_1, \dots, W_n) \sim \mathbf{P}$  ( $W_i$  corresponds to the representation of user  $i$ ) is obtained from website 1 and shared with colluding website 2.
3. Given the representation  $O$  of a random user on website 2, the attacker must find a  $W$ -measurable function to predict the identity of that representation to match that user on website 1.

**Modeling the distribution  $\mathcal{D}$**  A natural question is: how can we model the distribution  $\mathcal{D}$ ? And, can we define the optimal action  $\varphi$  taken by an attacker given their knowledge of  $\mathcal{D}$ ? Below we propose a natural parametric family of distributions  $\mathbb{D}$ , and an efficient way for the attacker to estimate their parameters based only on the samples known. We also describe a simple to implement optimal attack for this parametric family. In section 8.4 we verify empirically that our assumptions on  $\mathbb{D}$  closely match the observations from current web traffic.

Note that since  $\mathbf{P}$  is fully determined by the top sets  $S_i^s$ , we can equivalently define  $\mathbb{D}$  as a family of joint distributions over the sequence of random variables  $(S_i^s)_{i \in [n], s \in [r]}$ . We denote by  $\mathbb{P}$  the probability measure induced by  $\mathcal{D}$  on  $S_i^s$ . A distribution  $\mathcal{D}$  belongs to  $\mathbb{D}$  if it satisfies the following conditions:

1. All users have the same distribution for variable  $S_i^s$  for a given time  $s$  (i.e. the top topics of the user are sampled independent and identically distributed from the same distribution, but naturally users can have different sampled top topics.)

- For every user, samples of top sets are independent across time (but not necessarily identically distributed).

It is not hard to see that the assumptions on  $\mathbb{D}$  imply that distribution  $\mathcal{D} \in \mathbb{D}$  if and only if there exist distributions  $\mathcal{D}_1, \dots, \mathcal{D}_r$  such that

$$\mathcal{D}(\mathbf{P}) = \prod_{s=1}^r \mathcal{D}_s(\mathbf{P}_s).$$

For each distribution  $\mathcal{D}_s$  we will be interested in the following

$$\mathcal{D}_s(\mathbf{P}_s[i, o] = q_{in}) = \mathbb{P}(o \in S_i^s) = p_s[o]$$

These parameters represent the probability that topic  $o$  is part of the top 5 topics of a user. The following lemma, which is proved in the Appendix, shows how we can use observations from user representations to estimate these terms.

**Lemma 10.** *Let  $W_1^s, \dots, W_n^s$  be a sample of topics on website 1, let  $N = |\mathcal{T}|$  and  $\delta > 0$ . Let also*

$$\widehat{p}_s[o] = \frac{1}{n} \sum_{i=1}^n \frac{\mathbb{I}_{W_i^s=o} - q_{out}}{q_{in} - q_{out}},$$

where  $\mathbb{I}_{W_i^s=o}$  is 1 iff  $W_i^s = o$ . With probability  $1 - \delta$  uniformly across all topics  $o \in \mathcal{T}$  the following inequality holds:

$$|\widehat{p}_s[o] - p_s[o]| \leq \frac{1}{q_{in} - q_{out}} \sqrt{\frac{\log(2N/\delta)}{2n}}. \quad (5)$$

Let us now proceed to identify the optimal attacker under the assumptions on  $\mathcal{D}$ . We begin by discussing a natural attacker to derive some intuition into the results of our main theorem of the section.

**Example 1** (Hamming attack). *Given a representation  $o = (o^1, \dots, o^r)$  and representations  $W = (W_i^s)_{i \in [n], s \in [r]}$  a naïve attacker will assign representation  $o$  to a user  $i$  such that  $|\{s : W_i^s \neq o^s\}|$  is minimal. That is, it will naturally try to find the user that minimizes the Hamming distance between  $W_i$  and  $o$ .*

While the Hamming distance attack matches the intuition that the most likely user that generates a representation is that with the largest overlap on the topic sequence of website 1, one should keep in mind that not all topics are the same. Indeed, a match on a very unpopular topic should be worth more than a match on a popular topic. It is not hard to see that the parameters  $p_s[o]$  are a proxy for the popularity of a topic across the population. Thus an optimal attacker should leverage this information. The following theorem, proved in the Appendix, formalizes this intuition.

**Theorem 1** (Asymmetric Weighted Hamming Distance Attack). *Given a representation  $o = (o^1, \dots, o^r)$  on website 2, and representations  $W = (W_i^s)_{i \in [n], s \in [r]}$  on website 1. An attacker that wants to maximize its accuracy under the partial information setting selects the identity of the user that minimizes the following asymmetric weighted Hamming distance*

$$- \sum_{s: W_i^s = o^s} \log \left( q_{out} + \frac{(q_{in} - q_{out})q_{in}p_s[o^s]}{q_{out} + (q_{in} - q_{out})p_s[o^s]} \right) - \sum_{s: W_i^s \neq o^s} \log (q_{out} + (q_{in} - q_{out})\mathbb{P}(o^s \in S_i^s | W_i^s \in S_i^s)).$$

Notice how the attack in Theorem 1 can be seen as an asymmetric variant of the simple Hamming distance attack described above. It is important to note that while there may likely be several heuristics for utilizing the Topics API signal for re-identification our framework has allowed us to derive — from first principles — a simple optimal algorithm under some basic assumptions. We expect that future work on understanding the privacy of the Topics API can be done by relaxing some of these assumptions.

## 8 Empirical analysis

In this section we give empirical evaluations of our model for re-identification risk on two real-world datasets: (1) A Google proprietary dataset containing de-identified user data from a simulation of the Chrome Topics API and (2) the public Million Song Dataset [9]. In order to foster the reproducibility of our results, we released our code open source.<sup>5</sup> This section proceeds as follows. In sections 8.1-8.2, we present empirical implementations of the re-identification attacks presented before and more advanced machine learning heuristics. Next, in section 8.3, we show our results for the re-identification attack for Topics API. In section 8.4 we use machine learning methods to validate the hypothesis made in section 7. To do so, we further extend and validate the analysis of mutual information of Topics published previously [22]. In section 8.5, we show our results for the re-identification attack for Million Song Dataset.

### 8.1 Re-identification task for Topics API in random-user model

In this section we present our re-id attack on the Topics API. As a first step, we describe the data used

**Chrome Topics API data** In our empirical analysis we simulate the output observed by an adtech from the Topics API for a set of users over a period of time. This is achieved using a Google proprietary dataset of de-identified user browsing histories. Starting from this dataset, we run the Topics API algorithm for a sample of such users and simulate the output observed for two sites by the adtech from the Topics API sampling distribution. Our observation period consists of 8 intervals of 7 days of traffic shifted by 3 days each – that is we consider intervals  $[1, 7], [4, 10], [7, 13], \dots$ . Each interval (hence-forth epoch) is used to establish the top  $k$  topics of the user using the API topic model.<sup>6</sup> We restrict our analysis to the set of users that are observed in every epoch.

To be consistent with the current Topics API specification [26], for users with fewer than  $k$  topics in an epoch, we pad the top topics with random topics. Moreover, we set  $k = 5$  for the number of topics and use  $p = 5\%$  for the probability of releasing a random topic instead of an organic one, as currently implemented in the Chrome browser.

**Methods** We now present the methodology used to establish the accuracy of empirical re-identification attacks on the Topics API. We focus on the random-user model and consider different attack algorithms presented in section 8.2. We extract a target dataset of 10 million users chosen uniformly at random that are used in the re-identification attack analysis as the set  $\mathcal{I}$  over which to re-identify the user and simulate the observation on website 1 of all their topics sequence for  $r$  epochs ( $r \in [1, 8]$ ). Then we repeat 10,000 times a uniformly random draw of a user from the set  $\mathcal{I}$ , generate the  $r$ -length sequence on website 2 for such user and verify if the attack algorithm matches it correctly to the sample in website 1.

### 8.2 Attack algorithms

We simulate the following three attacks methods: the **Unweighted Hamming** attack and the **Asymmetric Weighted Hamming** attack as well as the **Neural Network** attack. We now describe each method before presenting our results in section 8.3.

#### 8.2.1 Unweighted Hamming Attack

This method is an exact implementation of the simple attack presented in Example 1.

---

<sup>5</sup>The code is available at: [https://github.com/google-research/google-research/tree/master/re\\_identification\\_risk](https://github.com/google-research/google-research/tree/master/re_identification_risk)

<sup>6</sup>The real Topics API has disjoint epochs of 7 days, here we simulate training over overlapping periods because our analysis is limited for privacy reasons to 4 weeks of data. This allows us to simulate longer topics sequences. We do not observe a significant difference in the results.

### 8.2.2 Asymmetric Weighted Hamming Attack

This method is a simplified implementation of the attack given in Theorem 1 which is optimal under the assumptions described above.

In our experiments, we further make the approximation of assuming that  $\forall W_i^s \neq o_s, \mathbb{P}(o^s \in S_i^s | W_i^s \in S_i^s)$  is a function only depending on  $o^s$ . Then we show that  $\forall W_i^s \neq o_s, \mathbb{P}(o^s \in S_i^s | W_i^s \in S_i^s) = \frac{4p_s[o^s]}{5-p_s[o^s]}$  (see Lemma 11 for more details). This reduces the parameters of the model to be estimated to only  $p_s[o^s]$ . Moreover, given that we empirically observe  $p_s[o^s]$  to be very close in every period we further assume  $p_s[o^s] = p[o^s]$ .

### 8.2.3 Neural Network Attack

In addition to the previous methods, we also implement a heuristic attack method based on a deep neural network. Our method is general and works on an arbitrary embedding representation for the topics sequence of a user. In section 8.4.1 we show how we obtain such embedding from sequence to sequence models while in this section we focus on how to use any embedding for matching users.

We train a deep neural network which takes two sequence embeddings as inputs and outputs a similarity score in  $[0, 1]$ , indicating the predicted probability that the two sequences are from the same user. Our network structure is similar to that of the Grale infrastructure [27]. The detailed network structure is presented in Figure 4.

**Training process.** We sample 20 million random users (different from the target set  $\mathcal{I}$  used in the re-identification task). For each sampled user  $u$ , we simulate a pair of topics sequences where both sequences are from  $u$  and we regard it as an example from the class of correct matching. We also create 10 pairs of sequences where the first sequence of each pair is the sequence from  $u$  and the second sequence is from a random user  $v \neq u$ , and we regard each pair as an example from the class of incorrect matching. Each training example is a pair of sequences generated by the above procedure and each sequence is embedded using a sequence model. The training objective of the neural network is to minimize the binary cross entropy loss.

**Re-identification inference.** Finally, to match a user sample using the neural network, given a user sequence  $A$ , we enumerate every sequence  $B$  from the target dataset and feed  $(A, B)$  to the neural network. We choose the sequence  $B^*$  as the re-identification output where the sequence  $(A, B^*)$  maximizes the predicted probability given by the neural network.

## 8.3 Results for the re-identification attack

We report the estimation of the probability of each attack algorithm to correctly re-identify the random user over the 10 million target set of users. In our analysis we study  $r = 1, 2, 4, 6, 8$  epochs. The main result is shown in Figure 3.

As expected, we observe that as the number of epochs increases there is an increased probability of correctly matching the user across the two sites (as more information is available to the attacker). Notice also how the more sophisticated Neural Network attack outperforms all methods but has performance close to our implementation of the provably good attack given by Theorem 1. This is a further confirmation of the validity of the simplifying assumptions made in our theoretical study. On the other hand, the simple unweighted Hamming attack performs less well. We also observe the gap between algorithms increases with the number of epochs. This is expected as more epochs of observation allow the advanced neural network algorithm to learn more correlations across the data. Overall, we observe that even after 8 epochs, the probability of correct re-identification is below 3%.

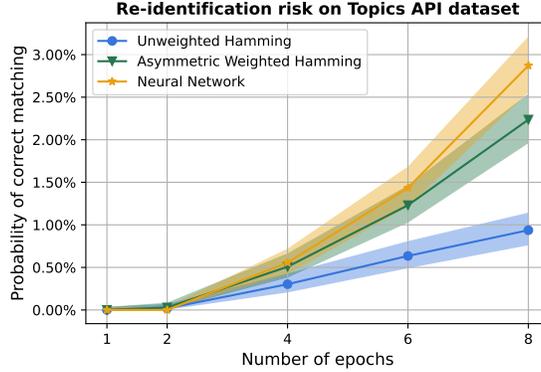


Figure 3: Probability of a correct cross-site match depending on the number of epochs observed. The 95% confidence intervals are reported. Notice how, even after 8 epochs, the probability of correct re-identification is below 3%.

### 8.4 Validation of the assumptions: Analysis of Mutual Information

In this section, our aim is to measure the validity of a key assumption made in Theorem 1 which we simulated: that of the cross-time independence of topics.

To do this we measure mutual information  $M^* = I(A^1, A^2, \dots, A^r; B^1, B^2, \dots, B^r)$ , for  $A = A^1, A^2, \dots, A^r$  and  $B = B^1, B^2, \dots, B^r$  being the two sequences of topics observed from a given user on two websites. Notice that, as shown in [22] in case of cross-time independence this  $M^*$  reduces to the easy to compute  $\sum_{s=1}^r I(A^s; B^s)$ . By estimating  $M^*$  and showing that it is close to  $\sum_{s=1}^r I(A^s; B^s)$  we verify the accuracy of the assumption.

Hyperparameter	Transformer	Transformer-LSTM
Encoder	Transformer	Transformer
Decoder	Transformer	LSTM
Attention dropout rate	0.1	0.1
Attention layer size	1,024	1,024
Dropout rate	0.1	0.1
Embedding size	1,024	1,024
MLP dimension	4,096	4,096
Number of attention heads	8	8
Number of encoder layers	6	6
Number of decoder layers	6	8
Decoder Hidden dimension	1,024	1,024
Training batch size	1,536	1,536
Total number of parameters	-	300M

Table 1: Hyper-parameters of the S2S Model Architectures.

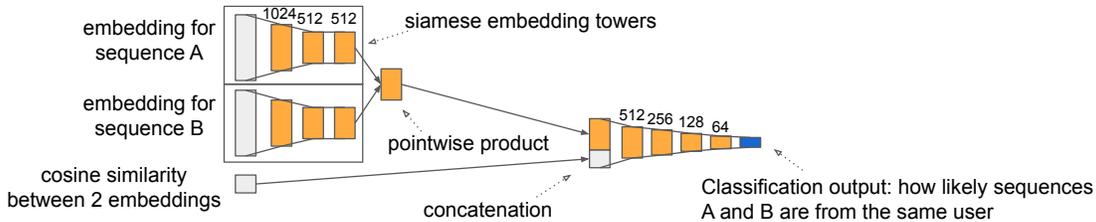


Figure 4: The structure of the network for re-identification attack. The embeddings of each sequence are computed by feeding the sequence to the trained S2S model (Transformer-Only) described in Section 8.4 and concatenating the hidden representations of the final encoder layer at each week. The number in the figure indicates the size of each layer. If not specified, each layer is a Fully Connected layer + ReLU activation.

In general estimation of  $M^*$  is a non-trivial task as the distribution of long sequences of topics has an exponentially growing support. In this section we use advanced ML techniques to tackle this challenge allowing us to extend the previously published work in [22] beyond the analysis of 2 epochs of data.

The rest of the section presents the ML model and how we used it to estimate the mutual information.

#### 8.4.1 Sequence to Sequence Model

We use state-of-the-art sequence to sequence models (S2S) [41] that are currently being used in a variety of machine learning applications ranging from natural language processing [43] to computer vision. As a byproduct of this model we develop an embedding for a topics sequence that we used in the previous section for the Neural Network attack.

First, we introduce formally the concept of sequence-to-sequence model. A sequence-to-sequence (S2S) model [41] assigns a probability to a sequence of target symbols  $B^1, B^2, \dots, B^Y$  given a sequence of source symbols  $A^1, A^2, \dots, A^X$ :  $P(B^1 = b^1, B^2 = b^2, \dots, B^Y = b^y | A^1 = a^1, A^2 = a^2, \dots, A^X = a^x)$ . The general architecture of an S2S model consists of an encoder network which generates an embedding of the source sequence and a decoder network that generates the target sequence conditional on the source sequence. A further refinement is an attention mechanism [6] that allows the decoder to attend to specific tokens of the source sequence when assigning probabilities to each token in the target sequence. A number of architectures have been proposed for the encoder and decoder networks. In this work, we will compare two popular architectures: Transformer [43] and a variant consisting of a Transformer encoder and a Long Short Term Memory (LSTM) decoder [12].

#### 8.4.2 Estimation method

Now, we present how to use S2S models to estimate mutual information.

The mutual information between sequences  $A$  and  $B$  can be written as a difference of two entropies [15]:

$$I(A^1, A^2, \dots, A^r; B^1, B^2, \dots, B^r) = H(B^1, B^2, \dots, B^r) - H(B^1, B^2, \dots, B^r | A^1, A^2, \dots, A^r) \quad (6)$$

Thus if we train two S2S models to estimate  $H(B^1, B^2, \dots, B^r)$  and  $H(A^1, A^2, \dots, A^r | B^1, B^2, \dots, B^r)$  respectively, we can compute an estimate of  $I(A^1, A^2, \dots, A^r; B^1, B^2, \dots, B^r)$ .

Given a training data set of tuples consisting of a user, and a pair of topics sequences associated with the user on two websites,  $A = A^1, A^2, \dots, A^r$ ,  $B = B^1, B^2, \dots, B^r$ , we estimate a sequence-to-sequence model to predict  $B$  given  $A$ . Over an unseen test data set consisting of  $n$  tuples  $\{A_i^1, A_i^2, \dots, A_i^r, B_i^1, B_i^2, \dots, B_i^r\}_{i=1}^n$ , we can estimate the conditional entropy of  $B$  given  $A$ :

$$H(B^1, B^2, \dots, B^r | A^1, A^2, \dots, A^r) = -\frac{1}{n} \sum_{i=1}^n \log P(B_i^1, B_i^2, \dots, B_i^r | A_i^1, A_i^2, \dots, A_i^r) \quad (7)$$

To estimate the unconditional entropy of the target sequence, we replace the source sequence within each tuple in the data set (both training and test) with a single special token \$, (i.e.  $A^1, A^2, \dots, A^r$  is replaced with \$). We can then use the above approach to train a model to estimate the unconditional entropy of the target sequence:

$$H(B^1, B^2, \dots, B^r) = -\frac{1}{n} \sum_{i=1}^n \log P(B_i^1, B_i^2, \dots, B_i^r | \$) \quad (8)$$

#### 8.4.3 Model Details

In our work we experiment with two separate S2S model architectures: either the vanilla Transformer [43] or a variant [12] consisting of a Transformer encoder and an LSTM decoder. (Exact details on the architecture

may be found in the above references). For estimating mutual information using a given sequence-to-sequence model architecture, we use the same set of hyper-parameters for both the unconditional and the conditional model (Table 1).

#### 8.4.4 Results

We now present the main set of results from the sequence-to-sequence model using the Transformer architecture. The hybrid architecture achieves a very close but slightly lower mutual information than the Transformer-only model so we omit its results and use the Transformer only embeddings in our re-identification analysis.

We measured the unconditional entropy of a sequence of topics for  $r$  epochs and as well as the conditional entropy (of the second sequence of the user) and their difference which is the mutual information. We report the results as bits/epoch. Using the Transformer model trained on  $r = 8$  epochs, we observe 6.54, 5.45 and 1.09 bits/epoch of unconditional entropy, conditional entropy and mutual information, respectively. On a single epoch of data we observe 7.59, 6.66, 0.93 bits for unconditional entropy, conditional entropy and mutual information, respectively.

Notice how, on average, we observe about 1.1 bits of mutual information per epoch of observation vs 0.93 bits of a single epoch. This suggests that while there is indeed some information gained by looking at sequences of topics across time, previous observations of the Topics API do not provide significant information about the topic returned in the current epoch. This validates our hypothesis that topics are close to independent across time and the model for the distribution  $\mathcal{D}$ .

### 8.5 Re-identification task for Million Song Dataset and Results

We now present our empirical study of re-identification attacks on the publicly available Million Song Dataset [9]. Similar to the study of Topics API, we focus on the random-user model.

In this dataset, a user is represented by all songs liked by them. The dataset contains 48 million entries on the listening activity of about a million users. The number of distinct songs is about a million.

We simulate a system that outputs a sample of  $r$  songs for a user, independently, to generate two different databases. Then, we measure the risk of re-identifying the user across the two datasets, depending on the number  $r$  of independent samples observed. The results are reported in Figure 5.

Here we use the Unweighted Hamming Attack for guessing the match of the user (notice that the other two attacks were specifically designed for Topics, so they are not meaningful here). We note that observing 4 independent random songs results in 1% re-identification risk. This is an example of how our framework allows us to assess the risk in such data releases.

## 9 Discussion and limitations

In this paper we have presented a framework for quantifying re-identification risk. Our theoretical formulation is general enough to frame many common notions of privacy like  $k$ -anonymity and differential privacy, and to capture real-world examples such as the Topics API. Our experiments show an empirical estimation of re-id risk even on very large representation spaces. We conclude this paper high-

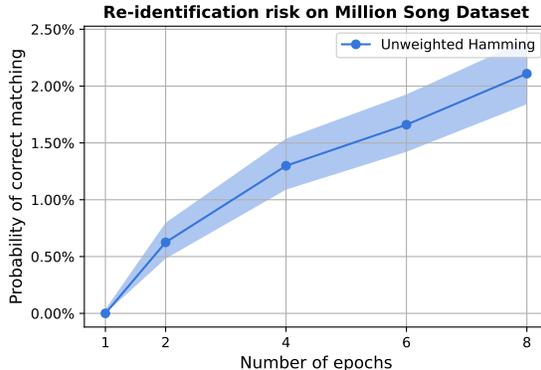


Figure 5: Probability of a correct cross-site match depending on the number of epochs observed. The 95% confidence intervals are reported.

lighting some of the limitations of our work and future directions of research in this area.

**Scope of the privacy risk measurement** First, it is important to observe that in our modeling we have focused exclusively, and purposefully, on re-identification risk. Preventing re-identification is only one of the many privacy and safety requirements of real-world applications. As real systems often require to bound additional risks (such as membership inference attacks [39]), we believe that privacy protections from differential privacy (DP) will be needed to complement re-identification risk analyses in many real systems.

**Theoretical limitations** From a theoretical point of view, we make mild modeling assumptions on the generation process for the representations. Specifically, we assume that each user sample from  $\mathbf{P}$  is drawn independently and we restrict our modeling to discrete representations. We remark that the partial information scenario we consider, and the distribution family  $\mathbb{D}$  we model, are just one of the many ways in which we could model an attacker with limited information on  $\mathbf{P}$ . We omit for instance, modeling the presence of more than 2 colluding clients.

More crucially, for the upper bound results, we make a closed-world assumption: we measure the re-identification risk under the assumption that this is the only information about a user available to an attacker. In other words, all of the side information is built into the representation matrix  $\mathbf{P}$ . While limiting in some situations, this assumption is reasonable in other cases (as we detail below). Moreover, our lower bounds hold even without this assumption (i.e., additional side information can only increase the risk).

**Limitations of the Topics API study** In generating the data for the Topics API study, we assumed that users visit both websites in all of the epochs; thus we specifically avoided dealing with gaps in the data. This is unlikely to be the case in a real world deployment of such an API. An additional concern is that we used overlapping time periods to simulate longer time horizons, to deal with limitations of the data we had on hand. On the modeling side, while we have tried to develop powerful re-identification methods, future researchers may be able to find more powerful approaches, increasing the empirical re-identification risk.

Finally, in our analysis we have focused exclusively on the Topics API output in isolation, ignoring other sources of side information that an attacker might have. Considering additional sources of information is beyond the scope of our work and would require incorporating this information into either the matrix  $\mathbf{P}$  or the distribution family  $\mathbb{D}$ .

**Applicability of the framework** In summary, our framework provides a lower bound on re-identification risk of data release that complements other privacy metrics that may be implemented by system designers, such as (differential) privacy budgets and aggregation (k-anonymity) guarantees.

In cases where the ‘closed-world’ assumption holds, we can use the methods developed in this work to provide even stronger upper bounds on re-identification risk. One potential such scenario of special interest is bounding *insider risk*: the risk of an employee maliciously re-identifying user information across a company’s systems. Given that such systems have strict controls on the flow of information (e.g. ACLs and auditing systems), we believe this framework can allow data protection officers to measure possibility of re-identification in such cases.

Overall, we expect that the tools we have developed could be used by privacy advocates for a data-driven verification of the privacy claims by data brokers and technology companies.

## A Proof of Theorem 1

*Proof.* By Lemma 2 we know that, given a representation  $o = (o_1, \dots, o_r)$ , the optimal attacker selects the user that maximizes

$$\max_i \mathbb{E}_{\mathcal{D}} [P[i, o] | W] = \max_i \prod \mathbb{E}_{\mathcal{D}_s} [P_s[i, o^s] | W_i^s], \quad (10)$$

where the last equality follows by the independence of representations across time and users. Let us calculate one of the terms in the above expression. Note that since  $P_s[i, o^s]$  can only take two values  $q_{\text{in}}$  and  $q_{\text{out}}$ , each factor in the above expectation is given by:

$$\begin{aligned} & \mathbb{E}_{\mathcal{D}_s} [P_s[i, o^s] | W_i^s] \\ &= \mathcal{D}_s(P_s[i, o^s] = q_{\text{in}} | W_i^s) q_{\text{in}} + \mathcal{D}_s(P_s[i, o^s] = q_{\text{out}} | W_i^s) q_{\text{out}} \\ &= q_{\text{out}} + \mathcal{D}_s(P_s[i, o^s] = q_{\text{in}} | W_i^s) (q_{\text{in}} - q_{\text{out}}) \end{aligned} \quad (11)$$

Note that if  $W_i^s = o^s$ , a straightforward application of Bayes rule and the fact that  $\mathbb{P}(W_i^s = o^s | P_s[i, o^s] = q_{\text{in}}) = q_{\text{in}}$  yields

$$\mathcal{D}_s(P_s[i, o^s] = q_{\text{in}} | W_i^s) = \frac{q_{\text{in}} p_s[o^s]}{q_{\text{in}} p_s[o^s] + q_{\text{out}} (1 - p_s[o^s])}$$

When  $W_i^s \neq o^s$  we simply rewrite the conditional expectation by the more explainable  $\mathbb{P}(o^s \in S_i^s | W_i^s \in S_i^s)$ . Plugging this expression in (11) we see that

$$\begin{aligned} E_{\mathcal{D}_s} [P_s[i, o^s] | W_i^s] &= \left( q_{\text{out}} + \frac{(q_{\text{in}} - q_{\text{out}}) q_{\text{in}} p_s[o^s]}{q_{\text{out}} + (q_{\text{in}} - q_{\text{out}}) p_s[o^s]} \right) \mathbb{I}_{W_i^s = o^s} \\ &\quad + (q_{\text{out}} + (q_{\text{in}} - q_{\text{out}}) \mathbb{P}(o^s \in S_i^s | W_i^s \in S_i^s)) \mathbb{I}_{W_i^s \neq o^s} \end{aligned}$$

The result follows by replacing this expression in (10) and taking the logarithm of the product.  $\square$

*Proof Of lemma 10.* Let  $\mathbf{W}$  be a random variable sampled as follows. Sample a set of top topics  $S \sim \mathcal{D}_s$ . Then sample  $\mathbf{W}$  according to

$$\mathbb{P}(\mathbf{W} = o | S) = \begin{cases} q_{\text{in}} & o \in S \\ q_{\text{out}} & o \notin S \end{cases}$$

It is very easy to see that the collection  $(W_i^s)$  is an i.i.d. sample from random variable  $\mathbf{W}$ . Moreover we have that

$$\begin{aligned} \mathbb{P}(\mathbf{W} = o) &= \sum_{S: o \in S} \mathbb{P}(\mathbf{W} = o | S) \mathcal{D}_s(S) + \sum_{S: o \notin S} \mathbb{P}(\mathbf{W} = o | S) \mathcal{D}_s(S) \\ &= q_{\text{in}} p_s[o] + (1 - p_s[o]) q_{\text{out}} \end{aligned}$$

Let  $Y_i[o] = \frac{\mathbb{I}_{W_i^s = o} - q_{\text{out}}}{q_{\text{in}} - q_{\text{out}}}$  and  $\hat{p}_s[o] = \frac{1}{n} \sum_{i=1}^n Y_i$ . Using the fact that  $|Y_i[o]| \leq \frac{1}{q_{\text{in}} - q_{\text{out}}}$  and  $\mathbb{E}[Y_i[o]] = p_s[o]$ , by Hoeffding's inequality we have with probability  $1 - \frac{\delta}{N}$ :

$$|p_s[o] - \hat{p}_s[o]| \leq \frac{1}{q_{\text{in}} - q_{\text{out}}} \sqrt{\frac{\log(2N/\delta)}{2n}}.$$

The result follows by using a union bound over all topics  $o$ .  $\square$

## B Modeling assumption on topics

In this section we discuss the modeling assumption on the Topics API used in the experiments (Section 8.1). The results of this section show that we can model an attacker using only the parameters  $p_s[o]$ .

**Assumption 1.** *For a fixed round  $s$  we assume that the random variable  $S_i^s$  representing the top set of users satisfies  $\mathbb{P}(o \in S_i^s | o' \in S_i^s) = \alpha_s(o)$  for all topics  $o' \neq o$ .*

The above assumption suggests some form of independence between the topics belonging to the top set. The following lemma shows that under this assumption  $\alpha_s(o)$  is in fact a simple function of  $p_s[o]$ .

**Lemma 11.** *Let  $S_i^s$  satisfy Assumption 1. Then*

$$\alpha_s(o) = \frac{4p_s[o]}{5 - p_s[o]}$$

*Proof.* By proposition 1 we know that:

$$\begin{aligned} 4p_s[o] &= \sum_{o' \neq o} \mathbb{P}((o, o') \in S_i^s) = \sum_{o' \neq o} \mathbb{P}(o \in S_i^s | o' \in S_i^s) P(o' \in S_i^s) \\ &= \alpha_s(o) \sum_{o' \neq o} p_s[o'] = \alpha_s(o)(5 - p_s[o]), \end{aligned}$$

where we have used Bayes rule and Assumption 1 for the second and third equalities respectively. The statement of lemma follows by rearranging terms.  $\square$

**Proposition 1.** *Let  $o$  be a fixed topic. The following properties holds for any distribution over the top set  $S_i^s$ .*

$$\sum_{o' \neq o} p_s[o'] = 5 - p_s[o] \tag{12}$$

$$\sum_{o' \neq o} \mathbb{P}((o, o') \in S_i^s) = 4p_s[o] \tag{13}$$

*Proof.* Fix a top set  $S_i^s$  of 5 elements. It is then easy to see that

$$\mathbb{I}_{o \in S_i^s} + \sum_{o' \neq o} \mathbb{I}_{o \in S_i^s} = 5 \tag{14}$$

Similarly, we claim that

$$\sum_{o' \neq o} \mathbb{I}_{o \in S_i^s} \mathbb{I}_{o' \in S_i^s} = 4\mathbb{I}_{o \in S_i^s} \tag{15}$$

Indeed, if  $o \notin S_i^s$  then the above expression is trivially true as both sides are 0. If, on the other hand,  $o \in S_i^s$ , then there are exactly 4 other elements in the set so  $\sum_{o' \neq o} \mathbb{I}_{o' \in S_i^s} = 4$ . The result follows by taking expectation of (15) and (13).  $\square$

## References

- [1] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pages 308–318. ACM, 2016.
- [2] Gagan Aggarwal, Tomas Feder, Krishnaram Kenthapadi, Rajeev Motwani, Rina Panigrahy, Dilys Thomas, and An Zhu. Approximation algorithms for k-anonymity. In *Proceedings of the International Conference on Database Theory (ICDT 2005)*, November 2005.
- [3] Rakesh Agrawal and Ramakrishnan Srikant. Privacy-preserving data mining. In *Proceedings of the 2000 ACM SIGMOD international conference on Management of data*, pages 439–450, 2000.
- [4] Mário S. Alvim, Konstantinos Chatzikokolakis, Annabelle McIver, Carroll Morgan, Catuscia Palamidessi, and Geoffrey Smith. *The Science of Quantitative Information Flow*. Information Security and Cryptography. Springer, 2020.

- [5] Ali Ismail Awad. Machine learning techniques for fingerprint identification: A short review. In Aboul Ella Hassanien, Abdel-Badeeh M. Salem, Rabie A. Ramadan, and Tai-Hoon Kim, editors, *Advanced Machine Learning Technologies and Applications - First International Conference, AMLTA 2012, Cairo, Egypt, December 8-10, 2012. Proceedings*, volume 322 of *Communications in Computer and Information Science*, pages 524–531. Springer, 2012.
- [6] Dzmitry Bahdanau, Kyung Hyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In *ICLR*, January 2015.
- [7] Muhammad Ahmad Bashir and Christo Wilson. Diffusion of user tracking data in the online advertising ecosystem. *Proc. Priv. Enhancing Technol.*, 2018(4):85–103, 2018.
- [8] Yoshua Bengio, Aaron C. Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(8):1798–1828, 2013.
- [9] Thierry Bertin-Mahieux, Daniel P.W. Ellis, Brian Whitman, and Paul Lamere. The million song dataset. In *12th International Society for Music Information Retrieval Conference (ISMIR 2011)*, 2011.
- [10] Richard Blahut. Hypothesis testing and information theory. *IEEE Transactions on Information Theory*, 20(4):405–417, 1974.
- [11] Kamalika Chaudhuri, Claire Monteleoni, and Anand D Sarwate. Differentially private empirical risk minimization. *Journal of Machine Learning Research*, 12(3), 2011.
- [12] Mia Xu Chen, Orhan Firat, Ankur Bapna, Melvin Johnson, Wolfgang Macherey, George Foster, Llion Jones, Mike Schuster, Noam Shazeer, Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Zhifeng Chen, Yonghui Wu, and Macduff Hughes. The best of both worlds: Combining recent advances in neural machine translation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 76–86, Melbourne, Australia, July 2018. Association for Computational Linguistics.
- [13] Aloni Cohen and Kobbi Nissim. Towards formalizing the gdpr’s notion of singling out. *Proceedings of the National Academy of Sciences*, 117(15):8344–8352, 2020.
- [14] Graham Cormode, Somesh Jha, Tejas Kulkarni, Ninghui Li, Divesh Srivastava, and Tianhao Wang. Privacy at scale: Local differential privacy in practice. In *Proceedings of the 2018 International Conference on Management of Data*, pages 1655–1658, 2018.
- [15] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. Wiley, 2001.
- [16] Mario Diaz, Hao Wang, Flavio P. Calmon, and Lalitha Sankar. On the robustness of information-theoretic privacy measures and mechanisms. *IEEE Transactions on Information Theory*, 66(4):1949–1978, 2020.
- [17] Cynthia Dwork. Differential privacy: A survey of results. In Manindra Agrawal, Ding-Zhu Du, Zhenhua Duan, and Angsheng Li, editors, *Theory and Applications of Models of Computation, 5th International Conference, TAMC 2008, Xi’an, China, April 25-29, 2008. Proceedings*, volume 4978 of *Lecture Notes in Computer Science*, pages 1–19. Springer, 2008.
- [18] Cynthia Dwork. Differential privacy and the us census. In *Proceedings of the 38th ACM SIGMOD-SIGACT-SIGAI symposium on principles of database systems*, pages 1–1, 2019.
- [19] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, pages 265–284. Springer, 2006.
- [20] Peter Eckersley. How unique is your web browser? In *Proceedings of the 10th International Conference on Privacy Enhancing Technologies*, page 1–18, Berlin, Heidelberg, 2010. Springer-Verlag.

- [21] Alessandro Epasto, Andrés Muñoz Medina, Steven Avery, Yijian Bai, Róbert Busa-Fekete, C. J. Carey, Ya Gao, David Guthrie, Subham Ghosh, James Ioannidis, Junyi Jiao, Jakub Lacki, Jason Lee, Arne Mauser, Brian Milch, Vahab S. Mirrokni, Deepak Ravichandran, Wei Shi, Max Spero, Yunting Sun, Umar Syed, Sergei Vassilvitskii, and Shuo Wang. Clustering for private interest-based advertising. In Feida Zhu, Beng Chin Ooi, and Chunyan Miao, editors, *KDD '21: The 27th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, Singapore, August 14-18, 2021*, pages 2802–2810. ACM, 2021.
- [22] Alessandro Epasto, Andres Munoz Medina, Christina Ilvento, and Josh Karlin. Measures of Cross-Site Re-Identification Risk: an Analysis of the Topics API Proposal. [https://github.com/patcg-individual-drafts/topics/blob/main/topics\\_analysis.pdf](https://github.com/patcg-individual-drafts/topics/blob/main/topics_analysis.pdf), 2022.
- [23] Alexandre Evfimievski, Johannes Gehrke, and Ramakrishnan Srikant. Limiting privacy breaches in privacy preserving data mining. In *Proceedings of the twenty-second ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 211–222, 2003.
- [24] Alejandro Gómez-Boix, Pierre Laperdrix, and Benoit Baudry. Hiding in the crowd: an analysis of the effectiveness of browser fingerprinting at large scale. In Pierre-Antoine Champin, Fabien Gandon, Mounia Lalmas, and Panagiotis G. Ipeirotis, editors, *Proceedings of the 2018 World Wide Web Conference on World Wide Web, WWW 2018, Lyon, France, April 23-27, 2018*, pages 309–318. ACM, 2018.
- [25] Google. The Privacy Sandbox. Available at [https://privacysandbox.com/intl/en\\_us/](https://privacysandbox.com/intl/en_us/) (Accessed October 15, 2022), 2022.
- [26] Google. The Topics API. Available at <https://github.com/patcg-individual-drafts/topics> (Accessed October 15, 2022), 2022.
- [27] Jonathan Halcrow, Alexandru Mosoi, Sam Ruth, and Bryan Perozzi. Grale: Designing networks for graph learning. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2523–2532, 2020.
- [28] Hsiang Hsu, Natalia Martinez, Martin Bertran, Guillermo Sapiro, and Flavio P. Calmon. A survey on statistical, information, and estimation—theoretic views on privacy. *IEEE BITS the Information Theory Magazine*, 1(1):45–56, 2021.
- [29] Shouling Ji, Prateek Mittal, and Raheem Beyah. Graph data anonymization, de-anonymization attacks, and de-anonymizability quantification: A survey. *IEEE Communications Surveys & Tutorials*, 19(2):1305–1326, 2016.
- [30] Batya Kenig and Tamir Tassa. A practical approximation algorithm for optimal k-anonymity. *Data Mining and Knowledge Discovery*, 25(1):134–168, 2012.
- [31] Pierre Laperdrix, Nataliia Bielova, Benoit Baudry, and Gildas Avoine. Browser fingerprinting: A survey. *ACM Transactions on the Web*, 14(2), 2020.
- [32] Kristen LeFevre, David J DeWitt, and Raghu Ramakrishnan. Mondrian multidimensional k-anonymity. In *22nd International conference on data engineering (ICDE'06)*, pages 25–25. IEEE, 2006.
- [33] Ashwin Machanavajjhala, Xi He, and Michael Hay. Differential privacy in the wild: A tutorial on current practices & open challenges. In *Proceedings of the 2017 ACM International Conference on Management of Data*, pages 1727–1730, 2017.
- [34] Frank McSherry and Kunal Talwar. Mechanism design via differential privacy. In *48th Annual IEEE Symposium on Foundations of Computer Science (FOCS'07)*, pages 94–103. IEEE, 2007.

- [35] Adam Meyerson and Ryan Williams. On the complexity of optimal k-anonymity. In *Proceedings of the twenty-third ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 223–228, 2004.
- [36] Lukasz Olejnik, Gunes Acar, Claude Castelluccia, and Claudia Diaz. The leaking battery. In *Data Privacy Management, and Security Assurance*, pages 254–263. Springer, 2015.
- [37] Hyounghmin Park and Kyuseok Shim. Approximate algorithms with generalizing attribute values for k-anonymity. *Information Systems*, 35(8):933–955, 2010.
- [38] Amrita Roy Chowdhury, Chenghong Wang, Xi He, Ashwin Machanavajjhala, and Somesh Jha. Crypte: Crypto-assisted differential privacy on untrusted servers. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*, pages 603–619, 2020.
- [39] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In *2017 IEEE symposium on security and privacy (SP)*, pages 3–18. IEEE, 2017.
- [40] Geoffrey Smith. On the foundations of quantitative information flow. In *International Conference on Foundations of Software Science and Computational Structures*, pages 288–302. Springer, 2009.
- [41] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112, 2014.
- [42] Latanya Sweeney. k-anonymity: A model for protecting privacy. *Int. J. Uncertain. Fuzziness Knowl. Based Syst.*, 10(5):557–570, 2002.
- [43] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5998–6008. Curran Associates, Inc., 2017.
- [44] Sergio Verdú.  $\alpha$ -mutual information. In *2015 Information Theory and Applications Workshop (ITA)*, pages 1–6, 2015.
- [45] Vassilios S Verykios, Elisa Bertino, Igor Nai Fovino, Loredana Parasiliti Provenza, Yucel Saygin, and Yannis Theodoridis. State-of-the-art in privacy preserving data mining. *ACM Sigmod Record*, 33(1):50–57, 2004.
- [46] Weina Wang, Lei Ying, and Junshan Zhang. On the relation between identifiability, differential privacy, and mutual-information privacy. *IEEE Transactions on Information Theory*, 62(9):5018–5029, 2016.
- [47] Genqiang Wu, Xianyao Xia, and Yeping He. Extending differential privacy for treating dependent records via information theory. *CoRR*, abs/1703.07474, 2017.
- [48] Mengmeng Yang, Lingjuan Lyu, Jun Zhao, Tianqing Zhu, and Kwok-Yan Lam. Local differential privacy and its applications: A comprehensive survey. *CoRR*, abs/2008.03686, 2020.