

LDPGuard: Defenses Against Data Poisoning Attacks to Local Differential Privacy Protocols

Kai Huang , Gaoya Ouyang , Qingqing Ye , Haibo Hu , Senior Member, IEEE, Bolong Zheng , Xi Zhao, Ruiyuan Zhang , and Xiaofang Zhou , Fellow, IEEE

Abstract—The protocols that satisfy Local Differential Privacy (LDP) enable untrusted third parties to collect aggregate information about a population without disclosing each user's privacy. In particular, each user locally encodes and perturbs his private data before sending it to the data collector, who aggregates and estimates the statistics about the population based on the collected perturbed values from individuals. Owing to their growing importance, LDP protocols have been widely studied and deployed in real-world scenarios (e.g., Chrome and Windows). However, as data poisoning attacks may be injected by attackers who introduce many fake users, the utility of the statistics is heavily poisoned. In this paper, we present a generic and extensible framework called LDPGuard to address the problem. LDPGuard provides effective defenses against data poisoning attacks to LDP protocols for frequency estimation, a basic query of most data analytics tasks. In particular, it first precisely estimates the percentage of fake users and then provides adversarial schemes to defend against particular data poisoning attacks. Experimental study on real-world and synthetic datasets demonstrates the superiority of LDPGuard compared to existing techniques.

Index Terms—Adversarial schemes, data poisoning attacks, frequency estimation, local differential privacy.

I. INTRODUCTION

LOCAL Differential Privacy (LDP) that enables an untrusted data collector to collect aggregate information about a population has been accepted as a *de facto* standard for data privacy in the local setting. In particular, each user locally encodes and perturbs his private data before sending it to the untrusted data collector. After receiving the distributed perturbed

data from users, the data collector designs detailed methods to aggregate and estimate statistics about the population without knowing the true values of each user. As the standard supports not only privacy-preserving data collection for untrusted third parties but also plausible deniability for each user, it has received a great amount of attention from both academia and industry. In recent years, a lot of protocols [4], [5], [6], [7], [8], [9], [10], [11], [13], [14], [15], [19] that satisfy Local Differential Privacy (a.k.a., LDP) have been developed. Some of the protocols have been deployed in industries such as Google, Apple, and Microsoft. In particular, Google integrated RAPPOR [7] into the Chrome browser to collect user responses to questions such as the default Chrome homepages while preserving their privacy. Apple [20] deployed LDP protocols in iOS to collect user preferences for emojis to identify popular emojis and recommend them to more users. Microsoft [23] developed LDP protocols to enable Windows 10 systems to aggregate statistics such as application usage. Most of these protocols are designed for frequency estimation, a basic query of most data analytics tasks (e.g., marginal release [5] and heavy hitter identification [11]). Since each user involved in the protocol is required to perturb his value before sending it to the data collector, it largely sacrifices the utility of analytics results obtained by the data collector. Therefore, almost all existing protocols focused on designing effective strategies to improve the utility of data.

However, these protocols are exposed to data poisoning attacks. Unlike poisoning attacks on machine learning models or systems [40], [41], [42], [43], [44], [45], [46], [47], [48], [49], [50], [65] where an attacker manipulates the training phase of a machine learning system by poisoning carefully selected training examples or tampering the training process, in poisoning attacks to LDP, an attacker introduces fake users to LDP protocols to craft the data and manipulate the data analytics results as they desire. In particular, fake users first choose some items (a.k.a., target items) and adopt various attack schemes to promote the estimated frequencies for target items. Previous studies [25], [61] have already discovered that attackers have access to a number of compromised accounts of web services (e.g., Hotmail, Twitter, and Google). In addition, they can purchase these compromised accounts from underground markets at very low prices (e.g., \$0.004 – \$0.03 for a Hotmail account and \$0.03 – \$0.50 for a phone verified Google account). In general, the mainstream poisoning attacks consist of random perturbed-value attack (RPA), random item attack (RIA) and maximal gain attack (MGA) [25]. For RPA, each fake user

Manuscript received 25 September 2022; revised 7 January 2024; accepted 23 January 2024. Date of publication 26 January 2024; date of current version 10 June 2024. This work was supported in part by the National Natural Science Foundation of China under Grants 62102334, 92270123, and 62372122, in part by the Research Grants Council, Hong Kong SAR, China under Grants 15226221, 15225921, 15209922, 15208923, and 15210023, and in part by General Research Grants under Grant FRG-24-027-FIE of the MUST Faculty Research Grants (FRG). Recommended for acceptance by X. Xiao. (*Corresponding author: Kai Huang.*)

Kai Huang is with the School of Computer Science and Engineering, Faculty of Innovation Engineering, Macau University of Science and Technology, Taipa, Macau (e-mail: kylehuangk@gmail.com).

Gaoya Ouyang, Xi Zhao, Ruiyuan Zhang, and Xiaofang Zhou are with the Department of Computer Science and Engineering, Hong Kong University of Science and Technology, Kowloon, Hong Kong (e-mail: ouyanggaoya@gmail.com; xzhaoca@connect.ust.hk; zry@ust.hk; zxf@cse.ust.hk).

Qingqing Ye and Haibo Hu are with the Department of Electronic and Information Engineering, Hong Kong Polytechnic University, Kowloon, Hong Kong (e-mail: qqing.ye@polyu.edu.hk; haibo.hu@polyu.edu.hk).

Bolong Zheng is with the School of Computer Science and Technology, Huazhong University of Science and Technology, Wuhan 430074, China.

Digital Object Identifier 10.1109/TKDE.2024.3358909

randomly selects a value from the encoded space of the LDP protocols and sends it to the data collector; for RIA, each fake user involved selects a value from the target items and perturbs it before sending it to the data collector; as for MGA, each fake user tries to maximize the frequency gains for the target items by crafting his value by solving an optimization problem (detailed in Section III). These attacks are posing serious security problems to existing LDP protocols and threats to LDP-based data analytics. For example, by applying data poisoning attacks to LDP protocols for default homepages in Chrome, attackers can promote a phishing webpage as the default homepage in Chrome. Similarly, they can degrade the user experience on emojis by increasing the frequencies of some emojis through data poisoning attacks. In addition, the popularity of malicious applications in Windows 10 systems could be increased in a similar way. In short, all these attacks can cause damage to user experience and commercial interests. In recent years, two main countermeasures [25] including *normalization* and *fake users detection* are proposed for defending against these attacks. The former is to normalize the estimated item frequencies such that each estimated item frequency is non-negative and the overall summation of item frequencies is 1, while the latter is to detect possible fake users based on their submitted values and to remove their impact on the final results. However, these countermeasures face two major challenges. The first is that the countermeasures are ad-hoc in nature, i.e., they can take effect on some scenarios but cannot adapt to others. The second is that they have limited effectiveness in reducing the impacts introduced by data poisoning attacks, as can be seen in Section VII.

In this paper, we present a unified framework called LDP-Guard to address the aforementioned challenges. In particular, LDPGuard adopts a relatively unified approach to handle different state-of-the-art LDP protocols for frequency estimation such as kRR [10], OUE [13], and OLH [13], and address all known data poisoning attacks. It also allows customization to address specific attacks and LDP protocols. It can be tailored according to the specific requirements of different attacks and protocols. LDPGuard primarily offsets attacks by observing statistical differences between two data collection instances. Therefore, this framework can also be applied to other potential attack methods and LDP protocols. Furthermore, LDPGuard adopts a two-round collection strategy to precisely estimate the percentage of fake users to facilitate effective countermeasures. It can also detect the attack type.

The rest of the paper is organized as follows. Related research is discussed in Section II, followed by preliminaries in Section III. The LDPGuard framework is described in Section IV and the estimation of the percentage of fake users is detailed in Section V. We introduce a detection method of poisoning attacks in Section VI. Experimental results are presented in Section VII. Section VIII concludes the paper. Formal proofs of all theorems and lemmas are given in [1].

II. RELATED WORK

Local Differential Privacy: Differential privacy (a.k.a., DP) is a standard that provides semantic, information-theoretic privacy guarantees, which is more stringent than traditional privacy techniques such as generalization [26], [27], [28], [29], [32].

DP was first proposed in the centralized setting with a trusted data collector [30]. In real-world life, users do not fully trust the data curator for privacy. Therefore, much attention has shifted to the local differential privacy (a.k.a., LDP), which can be cast as DP in the local setting and enables an untrusted data curator to collect aggregate information. A long line of research has been introduced for LDP especially frequency/mean estimation and practical data analysis task [2], [3], [4], [5], [6], [7], [8], [9], [10], [11], [12], [13], [14], [15], [16], [17], [18], [19], [21], [22], [31], [66], [67], [68], [69], [70], [71]. The state-of-the-art LDP protocols for frequency estimation consist of kRR [10], OUE [13], and OLH [13]. These protocols focused on improving data utility but may be exposed to data poisoning attacks, which are injected by attackers who introduce fake users to LDP protocols to craft the data and manipulate the analytics results as they desire.

Data Poisoning Attacks to LDP Protocols: Data poisoning attacks to LDP for frequencies consist of targeted attacks [25] and untargeted attacks [39]. The targeted attack, including random perturbed-value attack (RPA), random item attack (RIA), and maximal gain attack (MGA), increases estimated frequencies for the attacker-chosen items (i.e., targets). These attacks are posing serious security problems to existing LDP protocols and threats to LDP-based data analytics. In contrast, the untargeted attack distorts the item frequency distribution by manipulating the L_p -norm distance between frequency vectors before and after attacking. Recently, other data poisoning attacks have been developed, but they either focus on other data types (e.g., key-value data [73]) or do not specifically target frequency estimation (e.g., [72]).

Countermeasures Against Data Poisoning Attacks to LDP Protocols: Although there are some countermeasures against data poisoning attacks to LDP protocols [43], [51] and Sybil attacks [52], [53], [54], [55], [56], [57], [58], [59], [74], the former does not focus on LDP protocols for frequency estimation but machine learning models, the latter utilizes various information such as content, behavior, and social graphs to detect fake users in social networks. Recently, Cao et al. [25] present three countermeasures against data poisoning attacks to LDP protocols, namely, *normalization*, *fake users detection*, and *conditional probability based detection*. *Normalization* is to normalize the estimated item frequencies such that each estimated item frequency is non-negative and the summation of overall item frequencies is 1. *Fake users detection* is to detect potential fake users based on their submitted values and remove their impact on the final results. *Conditional probability based detection* is to detect possible fake users based on conditional probability, but it is applicable when there is only one target item. Moreover, the former two countermeasures are facing two major challenges. First, they are ad-hoc in nature, i.e., they can take effect on some scenarios but cannot adapt to others. Second, they have limited effectiveness in reducing the impacts introduced by data poisoning attacks.

III. PRELIMINARIES

Table I lists key notations and acronyms used in this paper. Without loss of generality, we assume that there are N users, each of which has a value $v \in \{1, 2, 3, \dots, d\}$ where

TABLE I
LIST OF KEY NOTATIONS

Notation	Description
$\epsilon, \epsilon_1, \epsilon_2$	privacy budgets
$d, \{1, 2, \dots, d\}$	number of items, original data space
\mathcal{D}	encoded data space
Encode(\cdot), Perturb(\cdot)	Encode algorithm, Decode algorithm
PE(\cdot)	composition operation of Encode and Perturb
SUP(y)	support set of the reported value y
p	$Pr[PE(v) \in \{y v \in \text{SUP}(y)\}]$
q	$Pr[PE(v') \in \{y v \in \text{SUP}(y)\}]$
p_1, q_1	p and q in the first round
p_2, q_2	p and q in the second round
p', q'	$p' = p_1$ or $p_2, q' = q_1$ or q_2
$\mathbb{I}_{\text{SUP}(y_i)}(v)$	indicator function
N, M	number of genuine users and fake users
$Y_{\text{true}}, Y_{\text{fake}}$	values from genuine users and fake users
y_i, Y	reported values from i -th user and all users
H, \mathcal{H}	a hash function, hash function family
r, τ	number of target items/ chosen target items
$T = \{t_1, t_2, \dots, t_r\}$	target items
$\tilde{f}_{t,b}, \tilde{f}_{t,a}$	estimated frequencies before/after attacking
f_t	frequency of target t over all genuine users
$\beta, \tilde{\beta}$	percentage/estimated percentage of fake users

$\{1, 2, \dots, d\}$ is **original data space** and d is the **number of items**. This paper focuses on LDP protocols for frequency estimation. The data collector aims to aggregate and estimate frequencies of values among the N users. To this end, the following three algorithms are performed one by one:

- **Encode**: A user who holds a value $v \in \{1, 2, \dots, d\}$ first encodes v to $\text{Encode}(v) \in \mathcal{D}$ where \mathcal{D} is **encoded data space** (i.e., the space of encoded values).
- **Perturb**: The user then perturbs the encoded value $\text{Encode}(v)$ to $\text{Perturb}(\text{Encode}(v))$ and reports it to the data collector. We can also use $\text{PE}(\cdot)$ to denote the composition operation of the Encode and Perturb algorithms for simplification.
- **Aggregate**: Data collector receives reported values from users and estimates the statistics of interest.

The Perturb algorithm should satisfy the following ϵ -local differential privacy (ϵ -LDP).

Definition 1. (Local Differential Privacy, LDP) A algorithm \mathcal{A} satisfies ϵ -LDP iff for any two inputs v and v' ,

$$\forall y \in \text{Range}(\mathcal{A}) : \frac{Pr[\mathcal{A}(v) = y]}{Pr[\mathcal{A}(v') = y]} \leq e^\epsilon \quad (1)$$

where $\epsilon (\geq 0)$ is the privacy budget and $\text{Range}(\mathcal{A})$ denotes the set of all possible outputs of the algorithm \mathcal{A} .

To analyze the accuracy of LDP protocols, the **Pure LDP** that supports simple and fast aggregation is developed [13].

Definition 2. (Pure Local Differential Privacy (Pure LDP)) A protocol \mathcal{A} with an Encode and Perturb function (i.e., $\text{PE}(\cdot)$) satisfies Pure Local Differential Privacy (or ϵ -Pure LDP) iff for any input value v a user holds,

$$\forall y \in \text{Range}(\mathcal{A}) : \frac{p}{q} \leq e^\epsilon$$

such that

$$\begin{aligned} Pr[PE(v) \in \{y|v \in \text{SUP}(y)\}] &= p \\ \forall v' \neq v : Pr[PE(v') \in \{y|v \in \text{SUP}(y)\}] &= q \end{aligned} \quad (2)$$

where $\text{SUP}(y)$ is the set of values that y supports [13], $\epsilon (\geq 0)$ is the privacy budget and $\text{Range}(\mathcal{A})$ is the set of all possible outputs of the protocol \mathcal{A} .

In particular, $\{y|v \in \text{SUP}(y)\}$ contains all outputs $\{y\}$ that “support” the occurrences of v (the support set of v for short). Note that the definition of the support set depends on a particular LDP protocol. For example, for the basic RAPPOR protocol [7] that encodes a value to a binary vector B whose v -th bit is 1, $\text{SUP}(B) = \{v|B[v] = 1\}$.

When each user follows a pure LDP to report his value y_i ($i \in \{1, 2, \dots, d\}$) to the data collector, the estimated frequency of v is

$$\tilde{f}_v = \frac{\frac{1}{N} \sum_{i=1}^N \mathbb{I}_{\text{SUP}(y_i)}(v) - q}{p - q} \quad (3)$$

where $\mathbb{I}_{\text{SUP}(y_i)}(v)$ is an indicator function. If $v \in \text{SUP}(y_i)$, $\mathbb{I}_{\text{SUP}(y_i)}(v) = 1$, and 0 otherwise. The estimated frequency \tilde{f}_v is an unbiased estimator of true frequency of v (Lemma 1), thus, $\mathbb{E}[\tilde{f}_v] = f_v$ and

$$\sum_{i=1}^N \mathbb{E}[\mathbb{I}_{\text{SUP}(y_i)}(v)] = N(f_v(p - q) + q). \quad (4)$$

Lemma 1: $\tilde{f}_v = \frac{\frac{1}{N} \sum_{i=1}^N \mathbb{I}_{\text{SUP}(y_i)}(v) - q}{p - q}$ is an unbiased estimator of the true frequency f_v of the item v .

A. Local Differential Privacy Protocols

In this subsection, we present three state-of-the-art LDP protocols for frequency estimation, namely kRR [10], OUE [13], and OLH [13].

1) **Randomized Response (kRR)**. **Encode**: Each user encodes his value v to itself, i.e., $\text{Encode}(v) = v$. The encoded data space $\mathcal{D} = \{1, 2, \dots, d\}$, which is the same as the original data space.

Perturb: The user keeps the true value v with probability p and perturbs it to a different value $v' \in \mathcal{D}$ with probability q , then reports the perturbed value $y = \text{PE}(v)$ to the data collector according to the following probability distribution.

$$Pr[y = a] = \begin{cases} \frac{e^\epsilon}{e^\epsilon + d - 1} \triangleq p & \text{if } a = v \\ \frac{1}{e^\epsilon + d - 1} \triangleq q & \text{otherwise} \end{cases} \quad (5)$$

Aggregate: The data collector receives the reported values $\{y_i|i \in \{1, 2, \dots, d\}\}$ from users, and estimates the frequency of a particular item v with Equation (3). In particular, the support set $\text{SUP}(y_i)(v)$ is $\{y|y = v\}$.

2) **Optimized Unary Encoding (OUE)**. **Encode**: Each user encodes his value v to a d -dimensional binary vector B where v -th bit is 1, i.e., $B[v] = 1$ and $B[v'] = 0$ for $v' \neq v$. The encoded data space $\mathcal{D} = \{0, 1\}^d$ since each bit could be 0 or 1.

Perturb: The user keeps each bit with value 1 (resp. 0) in B with probability p (resp. $1 - q$) and flips it with probability $1 - p$ (resp. q), then reports the perturbed value $y = \text{PE}(v)$ (or equally,

$y = \text{Perturb}(\mathbf{B})$) to the data collector. Formally,

$$\Pr[y[\text{pos}] = 1] = \begin{cases} \frac{1}{2} \triangleq p & \text{if } \mathbf{B}[\text{pos}] = 1 \\ \frac{1}{e^\epsilon + 1} \triangleq q & \text{if } \mathbf{B}[\text{pos}] = 0 \end{cases} \quad (6)$$

Aggregate: The data collector receives the reported values $\{y_i | i \in \{1, 2, \dots, d\}\}$ from users, and estimates the frequency of a particular item v with Equation (3). The support set $\text{SUP}(y_i)(v)$ for OUE is $\{y | y[v] = 1\}$.

3) Optimized Local Hashing (OLH). Encode: Given a universal hash function family \mathcal{H} such that the range of $H \in \mathcal{H}$ belongs to $\{1, 2, \dots, d'\}$, each user randomly selects a hash function H from \mathcal{H} and encodes his value v to $\langle H, H(v) \rangle$. The encoded data space $\mathcal{D} = \{\langle H, h \rangle | H \in \mathcal{H}, h \in \{1, 2, \dots, d'\}\}$. Note that \mathcal{H} is usually generated by the hash algorithm called xxHash [60], which implements the hash function family with different random seeds. For a better estimation performance, we follow existing work [13], [25] to set $d' = e^\epsilon + 1$.

Perturb: The user keeps the $H(v)$ with probability p and perturbs it to a different value $H(v)' \in \{1, 2, \dots, d'\}$ with probability q , and then reports the perturbed value $y = \text{PE}(v) = \langle H, H(v)' \rangle$ to the data collector. Formally,

$$\Pr[y = \langle H, h \rangle] = \begin{cases} \frac{e^\epsilon}{e^\epsilon + d' - 1} \triangleq p^* & \text{if } h = H(v) \\ \frac{1}{e^\epsilon + d' - 1} \triangleq q^* & \text{if } h \neq H(v) \end{cases} \quad (7)$$

Aggregate: The data collector receives the reported values $\{y_i | i \in \{1, 2, \dots, d\}\}$ from users, and estimates the frequency of a particular item v by (3) with the support set $\text{SUP}(y_i)(v) = \{y | y = \langle H, h \rangle \text{ and } H(v) = h\}$. Note that p^* and q^* are for hash function H instead of the pair of $\langle H, H(\cdot) \rangle$. As each hash function is randomly selected from \mathcal{H} , the overall perturbation probability parameters are $p = p^* = \frac{e^\epsilon}{e^\epsilon + d' - 1}$ and $q = \frac{1}{d'} p^* + (1 - \frac{1}{d'}) q^* = \frac{1}{d'}$.

B. Threat Model and Data Poisoning Attacks

Data poisoning attack is to increase the estimated frequencies for the attacker-chosen items (i.e., targets). In particular, given a target set with r target items, $T = \{t_1, t_2, \dots, t_r\}$ where $t_j \in \{1, 2, \dots, d\}$, suppose the estimated frequencies of a particular item t before and after attacking are $\tilde{f}_{t,b}$ and $\tilde{f}_{t,a}$ respectively, data poisoning attack is to injects M fake users to increase the **overall frequency gain** over T , i.e., $\sum_{t \in T} \Delta \tilde{f}_t = \tilde{f}_{t,a} - \tilde{f}_{t,b}$ where $\Delta \tilde{f}_t$ is the **frequency gain** of a single target item t .

In general, it consists of **random perturbed-value attack (RPA)**, **random item attack (RIA)**, and **maximal gain attack (MGA)** [25]. For RPA, each fake user randomly selects a value from the encoded space of LDP protocols and sends it to the data collector; for RIA, each fake user involved selects a value from the target item set and perturbs it before sending it to a data collector; as for MGA, each fake user tries to maximize the frequency gains for the target items by crafting his value via solving the following optimization problem. Let reported values of fake users be Y_{fake} , the overall frequency gain over T be

$\text{Gain}(Y_{\text{fake}})$, to maximize the gain, we have

$$\begin{aligned} \text{Max}_{Y_{\text{fake}}} \text{Gain}(Y_{\text{fake}}) &\iff \text{Max}_{Y_{\text{fake}}} \sum_{t \in T} \mathbb{E}[\Delta \tilde{f}_t] \\ &\iff \text{Max}_{Y_{\text{fake}}} \sum_{t \in T} \mathbb{E}[\tilde{f}_{t,a} - \tilde{f}_{t,b}] \iff \text{Max}_{Y_{\text{fake}}} \sum_{t \in T} \mathbb{E}[\tilde{f}_{t,a}] - \mathbb{E}[\tilde{f}_{t,b}] \\ &\stackrel{\text{Equ(3)}}{\iff} \text{Max}_{Y_{\text{fake}}} \sum_{t \in T} \mathbb{E} \left[\frac{\sum_{i=1}^{N+M} \mathbb{I}_{\text{SUP}(y_i)}(t)}{N+M} - q \right] \\ &\quad - \mathbb{E} \left[\frac{\sum_{i=1}^N \mathbb{I}_{\text{SUP}(y_i)}(t)}{N} - q \right] \\ &\iff \text{Max}_{Y_{\text{fake}}} \sum_{t \in T} \frac{\sum_{i=N+1}^{N+M} \mathbb{E}[\mathbb{I}_{\text{SUP}(y_i)}(t)]}{(N+M)(p-q)} - \frac{M \sum_{i=1}^N \mathbb{E}[\mathbb{I}_{\text{SUP}(y_i)}(t)]}{N(N+M)(p-q)} \\ &\stackrel{\text{Equ(4)}}{\iff} \text{Max}_{Y_{\text{fake}}} \frac{\sum_{i=N+1}^{N+M} \sum_{t \in T} \mathbb{E}[\mathbb{I}_{\text{SUP}(y_i)}(t)]}{(N+M)(p-q)} \\ &\quad - \frac{\sum_{t \in T} M(f_t(p-q) + q)}{(N+M)(p-q)}, \end{aligned} \quad (8)$$

where f_t is the frequency of the target t over all genuine users. As $\frac{\sum_{t \in T} M(f_t(p-q) + q)}{(N+M)(p-q)}$ only depends on genuine users, MGA is to craft Y_{fake} such that $\frac{\sum_{i=N+1}^{N+M} \sum_{t \in T} \mathbb{E}[\mathbb{I}_{\text{SUP}(y_i)}(t)]}{(N+M)(p-q)}$ is maximized. To this end, different implementation strategies are adopted for kRR, OUE, and OLH as listed below.

MGA for kRR: For kRR, observe that $\sum_{t \in T} \mathbb{E}[\mathbb{I}_{\text{SUP}(y_i)}(t)] \leq 1$ and $\sum_{t \in T} \mathbb{E}[\mathbb{I}_{\text{SUP}(y_i)}(t)] = 1$ if and only if y_i belongs to T . Therefore, MGA for kRR is obtained by randomly selecting a target item for each fake user.

MGA for OUE: For OUE, observe that $\sum_{t \in T} \mathbb{E}[\mathbb{I}_{\text{SUP}(y_i)}(t)] \leq r$ and $\sum_{t \in T} \mathbb{E}[\mathbb{I}_{\text{SUP}(y_i)}(t)] = r$ if and only if $\forall t \in T, y_i[t] = 1$. Intuitively, implementing MGA for OUE only needs to craft each fake user's value with $y_i[t] = 1$ for each target t , but this will result in only r bits 1 in the encoded binary vector. In contrast, each genuine user may report y_i with $\lfloor p + (d-1)q \rfloor$ bits 1. As such, MGA for OUE reports the binary vector with $y_i[t] = 1$ for each target t and $l = \lfloor p + (d-1)q - r \rfloor$ bits 1 for non-target bits.

MGA for OLH: For OLH, observe that $\sum_{t \in T} \mathbb{E}[\mathbb{I}_{\text{SUP}(y_i)}(t)] \leq r$ and $\sum_{t \in T} \mathbb{E}[\mathbb{I}_{\text{SUP}(y_i)}(t)] = r$ if and only if the selected hash function hashes all targets to the same value. In our implementation, we follow [25] to sample 1,000 hash functions from xxHash and select the one that results in the most targets to the same value.

IV. LDPGUARD: THE FRAMEWORK

In this section, we overview our LDPGuard framework for effective defenses against data poisoning attacks to local differential privacy protocols. In particular, LDPGuard first adopts a two-round collection method to accurately estimate the percentage of fake users, it then utilizes the estimated percentage

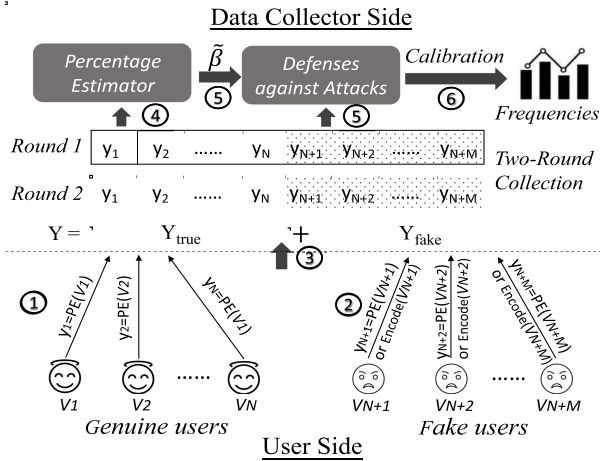


Fig. 1. LDPGuard framework.

to provide adversarial schemes to defend against particular data poisoning attacks such as RPA, RIA, and MGA.

A. Overview of LDPGuard

The workflow of LDPGuard is depicted in Fig. 1 and illustrated with Algorithm 1. Suppose there are N **genuine users**, each user holds a value v_i ($i \in \{1, 2, \dots, N\}$) in the **original data space** (i.e., $v_i \in \{1, 2, \dots, d\}$), encodes v_i to a value in **encoded data space** \mathcal{D} , perturbs the encoded value under ϵ_1 -LDP, and finally reports the perturbed value $y_i = PE(v_i)$ to the data collector (step ①, see Line 3, Algorithm 1). In addition, there are M **fake users**, each fake user holds a value v_i ($i \in \{N + 1, N + 2, \dots, N + M\}$) in the original data space, and adopts one of the following procedures to obtained reported value y_i (step ②, Line 5):

- *Encode only*: The fake user only encodes the input v_i to obtain y_i , i.e., $y_i = \text{Encode}(v_i)$. This applies to both RPA and MGA attacks, which do not utilize any LDP protocols to perturb user inputs. In particular, for RPA, each fake user randomly selects a value y_i from encoded data space and sends it to the data collector directly; for MGA, a fake user crafts a value y_i via solving the optimization problem (see (8)) and reports it without perturbation. In other words, each fake user with RPA and MGA can fabricate an input v_i and encode it to $y_i = \text{Encode}(v_i)$, which will be directly sent to the data collector without perturbation.
- *Encode and Perturb*: The fake user first encodes the input v_i and then perturbs the encoded value to obtain y_i , i.e., $y_i = \text{PE}(v_i)$. This applies to RIA attacks.

Therefore, the reported value y_i from a fake user is either $y_i = \text{PE}(v_i)$ or $y_i = \text{Encode}(v_i)$. The values $\{y_i | i \in \{1, 2, \dots, N + M\}\}$ are then reported to the data collector (step ③) for further processing (Lines 7 to 10).

Upon receiving the reported values Y consisting of $Y_{true} = \{y_i | i \in \{1, 2, \dots, N\}\}$ from genuine users and $Y_{fake} = \{y_i | i \in \{N+1, N+2, \dots, N+M\}\}$ from fake users in the first round, the second round with the same process is launched by the data collector (steps ①, ② and ③). Note that the privacy

Algorithm 1: LDPGuard.

Input: N inputs $\{v_1, v_2, \dots, v_N\}$ from genuine users and M fabricated values $\{v_{N+1}, v_{N+2}, \dots, v_{N+M}\}$ from fake users

Output: estimated frequencies of items

1: **for** $1 \leq round \leq 2$ **do**2: **for** $1 \leq i \leq N$ **do**3: $y_i = PE(v_i)$ 4: **for** $N + 1 \leq i \leq N + M$ **do**

5: $y_i = \text{Encode}(v_i)$ /*This is equivalent to directly selecting or fabricating a y_i from the encoded data space. */

6: $Y_{round} = \{y_i | 1 \leq i \leq N + M\}$

7: **if** attack is unknown **then**8: $attack \leftarrow \text{DETECTATTACK}(Y_1, Y_2)$

9: $\tilde{\beta} \leftarrow \text{PERCENTAGEESTIMATE}(Y_1, Y_2, \text{attack})$

$$10: \{\tilde{f}_v | v \in \{1, 2, \dots, d\}\} \leftarrow \text{DEFENSES}(Y_1, Y_2, \text{attack}, \tilde{\beta})$$
11: **return** $\{\tilde{f}_v | v \in \{1, 2, \dots, d\}\}$

budget in the second round is ϵ_2 and $\epsilon_1 + \epsilon_2 = \epsilon$ (by default, $\epsilon_1 = \epsilon_2 = \frac{1}{2}\epsilon$). We call this entire process “Two-Round Collection”. Then, the data collector utilizes **percentage estimator (detailed in Section V)** to estimate the percentage of fake users (step ④, Line 9) by comparing the reported values Y in the two rounds. The estimated percentage (denoted by $\tilde{\beta}$) is then used for **defenses against attacks (detailed in Section IV-B)** (step ⑤, Line 10). Finally, the data collector adopts calibrations (i.e., (3)) to derive item frequencies for publications (step ⑥, Line 10). The above process is based on the assumption that the data collector is aware of the attack. We shall remove the assumption and discuss the detection of attacks in Section VI (Line 8).

B. Defenses Against Data Poisoning Attacks

In this section, we assume the estimated percentage $\tilde{\beta}$ of fake users is already derived. We will remove this assumption in Section V. Based on $\tilde{\beta}$, LDPGuard provides effective adversarial schemes to defend against data poisoning attacks including RPA, RIA, and MGA to start-of-the-art LDP protocols such as kRR, OUE, and OLH. The main idea behind the adversarial schemes is to adopt the operation in the opposite direction of the one attacks adopted to offset the adverse effects of data poisoning attacks. In particular, LDPGuard first generates $(N + M)\tilde{\beta}$ records from possible distributions from which attacks sampled their values. It then removes the records from reported values of genuine and fake users to alleviate adverse effects of attacks. Detailed defenses for different attacks and LDP protocols are discussed below.

1) *Defenses Against Data Poisoning Attacks to kRR*: We begin by setting $p_1 = \frac{e^{\epsilon_1}}{e^{\epsilon_1} + d - 1}$, $q_1 = \frac{1}{e^{\epsilon_1} + d - 1}$, and $p_2 = \frac{e^{\epsilon_2}}{e^{\epsilon_2} + d - 1}$, $q_2 = \frac{1}{e^{\epsilon_2} + d - 1}$ where $\epsilon_1 + \epsilon_2 = \epsilon$. In addition, $p = \frac{e^{\epsilon}}{e^{\epsilon} + d - 1}$, $q = \frac{1}{e^{\epsilon} + d - 1}$.

Defenses Against RPA to kRR: First, LDPGuard randomly samples $(N + M)\tilde{\beta}$ values (denoted by \tilde{Y}_{fake}) with replacement from encoded data space $\mathcal{D} = \{1, 2, \dots, d\}$. For each sampled

value $v \in \tilde{Y}_{fake}$, LDPGuard then removes a corresponding record from all reported values Y if $v \in Y$. Finally, the data collector calibrates the results to obtain the estimated frequency for item v via

$$\tilde{f}_v = \frac{\frac{1}{|Y \setminus \tilde{Y}_{fake}|} \sum_{y_i \in Y \setminus \tilde{Y}_{fake}} \mathbb{I}_{\text{SUP}(y_i)}(v) - q'}{p' - q'} \quad (9)$$

where $p' = p_1$ or p_2 , $q' = q_1$ or q_2 , $Y \setminus \tilde{Y}_{fake}$ is the remaining records in Y after removing \tilde{Y}_{fake} . The expectation of \tilde{f}_v is

$$\begin{aligned} \mathbb{E}[\tilde{f}_v] &= \frac{\mathbb{E} \left[\frac{1}{|Y \setminus \tilde{Y}_{fake}|} \sum_{y_i \in Y \setminus \tilde{Y}_{fake}} \mathbb{I}_{\text{SUP}(y_i)}(v) \right] - q'}{p' - q'} \\ &= \frac{\frac{N_1(f_v(p' - q') + q') + N_2 f_v^*}{N_1 + N_2} - q'}{p' - q'} \end{aligned} \quad (10)$$

where $N_1 = |Y_{true} \cap (Y \setminus \tilde{Y}_{fake})|$, $N_2 = |Y \setminus \tilde{Y}_{fake} \setminus Y_{true}|$, and $f_v^* = \frac{1}{d}$ is the probability that v equals the value sampled from $Y \setminus \tilde{Y}_{fake} \setminus Y_{true}$. Observe that if $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{fake} = Y_{true}$), $N_2 = 0$, \tilde{f}_v is an unbiased estimator of f_v since $\mathbb{E}[\tilde{f}_v] = f_v$. The overall frequency gain is

$$\begin{aligned} \text{Gain}(Y_{fake}) &= \sum_{t \in T} \mathbb{E}[\tilde{f}_{t,a} - \tilde{f}_{t,b}] = \sum_{t \in T} \mathbb{E}[\tilde{f}_{t,a}] - \mathbb{E}[\tilde{f}_{t,b}] \\ &\stackrel{\text{Equ(3)}}{\stackrel{\text{Equ(10)}}{=}} \sum_{t \in T} \frac{\frac{N_1(f_t(p' - q') + q') + N_2 * 1/d}{N_1 + N_2} - q'}{p' - q'} \\ &\quad - \mathbb{E} \left[\frac{\sum_{i=1}^N \mathbb{I}_{\text{SUP}(y_i)}(t)}{N} - q \right] \\ &\stackrel{\text{Equ(4)}}{=} \sum_{t \in T} \frac{\frac{N_1(f_t(p' - q') + q') + N_2 * 1/d}{N_1 + N_2} - q'}{p' - q'} - \frac{(f_t(p - q) + q) - q}{p - q} \\ &= \sum_{t \in T} \frac{\frac{N_1(f_t(p' - q') + q') + N_2 * 1/d}{N_1 + N_2} - f_t(p' - q')}{p' - q'}. \end{aligned} \quad (11)$$

If $N_1 = N$, $N_2 = 0$, the gain $\text{Gain}(Y_{fake})$ is 0.

Defenses Against RIA to kRR: For RIA, LDPGuard first randomly samples $(N + M)\tilde{\beta}$ values (denoted by \tilde{Y}_{fake}) with replacement from target set $T = \{t_1, t_2, \dots, t_r\}$, and then encodes each value to the encoded data space $\mathcal{D} = \{1, 2, \dots, d\}$. For the encoded data, it keeps the true value with probability $p = p'$ ($p' = p_1$ or p_2) and perturbs it to a different value in \mathcal{D} with probability $q = q'$ ($q' = q_1$ or q_2). Finally, LDPGuard removes the perturbed value from Y if the perturbed value is found in Y_2 and calibrates the results to obtain the estimated frequency \tilde{f}_v of the item v using Equation (9). According to Equation (10), the expectation of \tilde{f}_v is

$$\begin{aligned} \mathbb{E}[\tilde{f}_v] &= \frac{\frac{N_1(f_v(p' - q') + q') + N_2 f_v^*}{N_1 + N_2} - q'}{p' - q'}, \\ f_v^* &= \begin{cases} \frac{1}{r}(p' - q') + q', & v \in T \\ 0, & v \notin T \end{cases} \end{aligned} \quad (12)$$

where $N_1 = |Y_{true} \cap (Y \setminus \tilde{Y}_{fake})|$, $N_2 = |Y \setminus \tilde{Y}_{fake} \setminus Y_{true}|$, and f_v^* is the probability that v equals the value sampled from $Y \setminus \tilde{Y}_{fake} \setminus Y_{true}$. If $N_1 = N$, $N_2 = 0$, \tilde{f}_v is an unbiased estimator of f_v since $\mathbb{E}[\tilde{f}_v] = f_v$, and overall frequency gain $\text{Gain}(Y_{fake})$ shown below is 0.

$$\begin{aligned} \text{Gain}(Y_{fake}) &= \sum_{t \in T} \mathbb{E}[\tilde{f}_{t,a} - \tilde{f}_{t,b}] = \sum_{t \in T} \mathbb{E}[\tilde{f}_{t,a}] - \mathbb{E}[\tilde{f}_{t,b}] \\ &\stackrel{\text{Equ(4)}}{\stackrel{\text{Equ(12)}}{=}} \sum_{t \in T} \frac{\frac{N_1(f_t(p' - q') + q') + N_2(\frac{p' - q'}{r} + q')}{N_1 + N_2} - q'}{p' - q'} \\ &\quad - \frac{f_t(p - q) + q - q}{p - q} \\ &= \sum_{t \in T} \frac{\frac{N_1(f_t(p' - q') + q') + N_2(1/r(p' - q') + q')}{N_1 + N_2} - f_t(p' - q')}{p' - q'}. \end{aligned}$$

Defenses Against MGA to kRR: For RPA, LDPGuard first randomly samples $(N + M)\tilde{\beta}$ values (denoted by \tilde{Y}_{fake}) with replacement from target set $T = \{t_1, t_2, \dots, t_r\}$ instead of \mathcal{D} . For each value sampled, LDPGuard then removes the record $v \in \tilde{Y}_{fake}$ from Y if $v \in Y$. Finally, the data collector calibrates the results to obtain the estimated frequency \tilde{f}_v for the item v using (9). According to Equation (10), the expectation of \tilde{f}_v is

$$\mathbb{E}[\tilde{f}_v] = \frac{\frac{N_1(f_v(p' - q') + q') + N_2 f_v^*}{N_1 + N_2} - q'}{p' - q'}, f_v^* = \begin{cases} \frac{1}{r} & \text{if } v \in T \\ 0 & \text{if } v \notin T \end{cases} \quad (13)$$

The overall frequency gain

$$\begin{aligned} \text{Gain}(Y_{fake}) &= \sum_{t \in T} \mathbb{E}[\tilde{f}_{t,a} - \tilde{f}_{t,b}] = \sum_{t \in T} \mathbb{E}[\tilde{f}_{t,a}] - \mathbb{E}[\tilde{f}_{t,b}] \\ &\stackrel{\text{Equ(4)}}{\stackrel{\text{Equ(13)}}{=}} \sum_{t \in T} \frac{\frac{N_1(f_t(p' - q') + q') + N_2 * 1/r}{N_1 + N_2} - q'}{p' - q'} \\ &\quad - \frac{(f_t(p - q) + q) - q}{p - q} \\ &= \sum_{t \in T} \frac{\frac{N_1(f_t(p' - q') + q') + N_2 * 1/r}{N_1 + N_2} - f_t(p' - q')}{p' - q'}. \end{aligned}$$

Accordingly, $\text{Gain}(Y_{fake})$ is 0 if $N_1 = N$ and $N_2 = 0$.

2) **Defenses Against Data Poisoning Attacks to OUE:** For OUE, we set $p_1 = \frac{1}{2}$, $q_1 = \frac{1}{e\epsilon_1 + 1}$, and $p_2 = \frac{1}{2}$, $q_2 = \frac{1}{e\epsilon_2 + 1}$ where $\epsilon_1 + \epsilon_2 = \epsilon$. Also, let $p = \frac{1}{2}$, $q = \frac{1}{e\epsilon + 1}$, $p' = p_1$ or p_2 , and $q' = q_1$ or q_2 .

Defenses Against RPA to OUE: LDPGuard first samples $(N + M)\tilde{\beta}$ values (denoted by \tilde{Y}_{fake}) randomly with replacement from the encoded data space $\mathcal{D} = \{0, 1\}^d$. For each value sampled, LDPGuard then removes the record $v \in \tilde{Y}_{fake}$ from all reported Y if $v \in Y$. Finally, the data collector calibrates the results to obtain the estimated frequency \tilde{f}_v for the item v using (9). On average, each value sampled from $\mathcal{D} = \{0, 1\}^d$ supports

v with probability $\frac{1}{2}$, according to (10), the expectation of \tilde{f}_v is

$$\mathbb{E}[\tilde{f}_v] = \frac{\frac{N_1(f_v(p'-q')+q')+N_2*1/2}{N_1+N_2} - q'}{p' - q'} \quad (14)$$

where $N_1 = |Y_{true} \cap (Y \setminus \tilde{Y}_{fake})|$ and $N_2 = |Y \setminus \tilde{Y}_{fake} \setminus Y_{true}|$. Observe that if $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{fake} = Y_{true}$), $N_2 = 0$, \tilde{f}_v is an unbiased estimate of f_v since $\mathbb{E}[\tilde{f}_v] = f_v$ and the following overall frequency gain $Gain(Y_{fake})$ is 0.

$$Gain(Y_{fake}) \stackrel{Equ(14)}{=} \sum_{t \in T} \frac{\frac{N_1(f_t(p'-q')+q')+N_2*1/2}{N_1+N_2} - f_t(p' - q')}{p' - q'}.$$

Defenses Against RIA to QUE: For RIA, LDPGuard first randomly samples $(N + M)\tilde{\beta}$ values (denoted by \tilde{Y}_{fake}) with replacement from target set $T = \{t_1, t_2, \dots, t_r\}$, and then encodes each value to the encoded data space $\mathcal{D} = \{0, 1\}^d$. Obviously, each encoded value is a d -bits binary vector, where each bit with value 1 (resp. 0) is kept with probability p' = p_1 or p_2 (resp. $1 - q'$) and flipped with probability $1 - p'$ (resp. $q' = q_1$ or q_2). Finally, LDPGuard removes the perturbed value from Y if the perturbed value is found in Y , and calibrates the results to obtain the estimated frequency for the item v using Equation (9). Similar to RIA to kRR, the expectation of \tilde{f}_v is $\frac{(N_1(f_v(p'-q')+q')+N_2*f_v^*)/(N_1+N_2)-q')}{p'-q'}$ where $f_v^* = (1/r(p' - q') + q')$, if $v \in T$; and $f_v^* = q'$, otherwise. Consequently, the overall frequency gain is

$$Gain(Y_{fake}) = \sum_{t \in T} \frac{\frac{N_1(f_t(p'-q')+q')+N_2(1/r(p'-q')+q')}{N_1+N_2} - f_t(p' - q')}{p' - q'}.$$

Similarly, $Gain(Y_{fake})$ is 0 if $N_1 = N$ and $N_2 = 0$.

Defenses against MGA to OUE: As for MGA to OUE, LDPGuard first initializes $(N + M)\tilde{\beta}$ d -bits zero vectors (denoted by \tilde{Y}_{fake}), for $y_i \in \tilde{Y}_{fake}$, and then sets $y_i[t] = 1$ for each bit $t \in T$ and $l = \lfloor p + (d - 1)q - r \rfloor$ randomly selected non-target bits. After that, LDPGuard removes the record $y_i \in \tilde{Y}_{fake}$ from Y if $y_i \in Y$. Finally, the data collector calibrates the results to obtain the estimated frequency of the item v using (9). The expectation of \tilde{f}_v is $\frac{(N_1(f_v(p'-q')+q')+N_2*f_v^*)/(N_1+N_2)-q')}{p'-q'}$ where $f_v^* = 1$, if $v \in T$; and $f_v^* = 0$, otherwise. The overall frequency gain is

$$Gain(Y_{fake}) = \sum_{t \in T} \frac{\frac{N_1(f_t(p'-q')+q')+N_2}{N_1+N_2} - f_t(p' - q')}{p' - q'}.$$

\tilde{f}_v is an unbiased estimator of f_v and the overall frequency gain $Gain(Y_{fake})$ is 0 if $N_1 = N$ and $N_2 = 0$.

3) Defenses Against Data Poisoning Attacks to OLH: Let $p^* = \frac{e^\epsilon}{e^\epsilon + d' - 1}$, $q^* = \frac{1}{e^\epsilon + d' - 1}$, $p = p^*$ and $q = \frac{1}{d'}p^* + (1 - \frac{1}{d'})q^* = \frac{1}{d'}$. Similarly, p_1^* , q_1^* , p_1 and q_1 (resp. p_2^* , q_2^* , p_2 and q_2) are defined w.r.t. ϵ_1 (resp. ϵ_2). In addition, let $p' = p_1$ or p_2 , $q' = q_1$ or q_2 .

Defenses Against RPA to OLH: First, LDPGuard crafts $(N + M)\tilde{\beta}$ key-value pairs $\{\langle H, h \rangle | H \in \mathcal{H}, h \in \{1, 2, \dots, d'\}\}$ (denoted by \tilde{Y}_{fake}) where H is randomly selected from \mathcal{H} and

h is randomly drawn from $\{1, 2, \dots, d'\}$ ($d' = e^\epsilon + 1$). Each record in \tilde{Y}_{fake} is removed from Y if it exists in Y . Finally, the data collector calibrates the results to obtain the frequency of the item v using (9). The expectation of \tilde{f}_v is $\frac{(N_1(f_v(p'-q')+q')+N_2*f_v^*)/(N_1+N_2)-q')}{p'-q'}$ where $f_v^* = \frac{1}{d'}$, if $v \in T$; and $f_v^* = 0$, otherwise. The overall frequency gain

$$Gain(Y_{fake}) = \sum_{t \in T} \frac{\frac{N_1(f_t(p'-q')+q')+N_2*1/d'}{N_1+N_2} - f_t(p' - q')}{p' - q'}.$$

If $N_1 = N$, $N_2 = 0$, \tilde{f}_v is an unbiased estimator of f_v since $\mathbb{E}[\tilde{f}_v] = f_v$. In addition, the overall frequency gain $Gain(Y_{fake})$ is 0.

Defenses Against RIA to OLH: To defend against RIA to OLH, LDPGuard first samples $(N + M)\tilde{\beta}$ values from target set \mathcal{T} . For each sampled value v , a hash function H is randomly selected from \mathcal{H} . Then, the key-value pair $\langle H, h \rangle$ is perturbed to $\langle H, h' \rangle$ ($h' \neq h$) with probability q_1^* or q_2^* , and remains unchanged with probability p_1^* or p_2^* . All these perturbed values are removed from Y , and the data collector calibrates the results to obtain the estimated frequency for the item v using (9). The expectation of \tilde{f}_v is $\frac{(N_1(f_v(p'-q')+q')+N_2*f_v^*)/(N_1+N_2)-q')}{p'-q'}$ where $f_v^* = \frac{1}{r}(p' - q') + q'$, if $v \in T$; and $f_v^* = q'$, otherwise. The overall frequency gain is

$$Gain(Y_{fake}) = \sum_{t \in T} \frac{\frac{N_1(f_t(p'-q')+q')+N_2*(\frac{1}{r}(p'-q')+q')}{N_1+N_2} - f_t(p' - q')}{p' - q'}.$$

Similarly, \tilde{f}_v is an unbiased estimator of f_v since $\mathbb{E}[\tilde{f}_v] = f_v$, and $Gain(Y_{fake})$ is 0 if $N_1 = N$ and $N_2 = 0$.

Defenses Against MGA to OLH: As for MGA to OLH, LDPGuard crafts $(N + M)\tilde{\beta}$ key-value pairs $\{\langle H, h \rangle | H \in \mathcal{H}, h \in \{1, 2, \dots, d'\}\}$ (denoted by \tilde{Y}_{fake}) where H is randomly selected from \mathcal{H} and h is the hash function that results in the most targets to the same value. Each record in \tilde{Y}_{fake} is removed from Y if it exists in Y . Finally, the data collector calibrates the results to obtain the estimated frequency for the item v with (9). The expectation of \tilde{f}_v is $\frac{(N_1(f_v(p'-q')+q')+N_2*f_v^*)/(N_1+N_2)-q')}{p'-q'}$ where $f_v^* = 1$, if $v \in T$; and $f_v^* = 0$, otherwise. Obviously, if $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{fake} = Y_{true}$), $N_2 = 0$, \tilde{f}_v is an unbiased estimator of f_v since $\mathbb{E}[\tilde{f}_v] = f_v$, and the following gain is 0.

$$Gain(Y_{fake}) = \sum_{t \in T} \frac{\frac{N_1(f_t(p'-q')+q')+N_2}{N_1+N_2} - f_t(p' - q')}{p' - q'}.$$

C. Discussions

Communication and Computing Cost: The communication cost of each user is $\mathcal{O}(\log d)$, $\mathcal{O}(d)$, and $\mathcal{O}(\log d)$ for kRR, OUE and OLH, respectively. The computing cost of the data collector is $\mathcal{O}(N + M + d)$, $\mathcal{O}((N + M) * d)$, and $\mathcal{O}((N + M) * d)$ for kRR, OUE and OLH, respectively.

Privacy Analysis: The first round collection satisfies ϵ_1 -LDP and the second round collection satisfies ϵ_2 -LDP, according to

Sequential Composition (see Lemma 2), LDPGuard satisfies ϵ -LDP where $\epsilon = \epsilon_1 + \epsilon_2$.

Lemma 2: Given C algorithms $\{\mathcal{A}_i | i \in \{1, 2, \dots, C\}\}$, each \mathcal{A}_i satisfies ϵ_i -LDP, the sequence of algorithms \mathcal{A}_i ($i \in \{1, 2, \dots, C\}$) collectively satisfies $\sum_{i \in \{1, 2, \dots, C\}} \epsilon_i$ -LDP.

Furthermore, observe that if $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{fake} = Y_{true}$), $N_2 = 0$, \tilde{f}_v is an unbiased estimator of f_v since $\mathbb{E}[\tilde{f}_v] = f_v$, and the overall frequency gain $Gain(Y_{fake})$ is 0. In other words, the effects of poisoning attacks are offset with the proposed defenses if $N_1 = N$ (i.e., $Y \setminus \tilde{Y}_{fake} = Y_{true}$). To this end, \tilde{Y}_{fake} should be sampled from the same distribution as the one injected into attacks, as mentioned in this section. Furthermore, the same number of noise data should be removed, i.e., $(N + M)\tilde{\beta}$ should be close to the real number of fake users (i.e., M). We can address this problem by accurately estimating the percentage of fake users as discussed in Section V.

V. ESTIMATION OF THE PERCENTAGE OF FAKE USERS

Observe that the reported values of fake users may be statistically abnormal compared to those of genuine users. For example, the values of fake users who engaged in RPA attacks are very randomly distributed in each round of the two-round collection; those of fake users with MGA attacks are too determinate so that they can be distinguished from the reported values from genuine users. This observation motivates us to accurately estimate the percentage of fake users by comparing the values reported in the two-round collection. In particular, let P_1 (resp. P_2) be the probability that the data collector collects the same value from the same genuine user (resp. fake user) in two different rounds, and CNT is the number of users who reported the same value in two different rounds, the percentage of fake users (i.e., $\beta = \frac{M}{N+M}$) can be estimated by

$$\tilde{\beta} = \frac{(N + M)P_1 - CNT}{(N + M)(P_1 - P_2)}. \quad (15)$$

Theorem 1: $\tilde{\beta} = \frac{(N+M)P_1 - CNT}{(N+M)(P_1 - P_2)}$ is an unbiased estimator of the percentage of fake users.

The problem now becomes how to compute P_1 and P_2 . We answer this question in the remaining parts of this section.

A. Percentage Estimation for kRR

Recall that each genuine user reports his value with kRR follows the following three steps: 1) encodes his value v to itself, i.e., $\text{Encode}(v) = v$; 2) keeps the true value v with probability p' and perturbs it to a different value $v' \in \mathcal{D} = \{1, 2, \dots, d\}$ with probability q' . For the two-round collection, let $p' = p_1$ or p_2 and $q' = q_1$ or q_2 . 3) and finally reports the perturbed value $y = \text{PE}(v)$ to the data collector. Therefore, a genuine user reports the same value when he 1) reports the true value v in both the first and second round; or 2) reports the same value v' ($v' \neq v$). The probability P_1 that the data collector collects the same value from a genuine user in two rounds is

$$P_1 = (p')^2 + (d - 1)(q')^2. \quad (16)$$

where $(p')^2$ (resp. $(q')^2$) is the probability that a genuine user reports the same value v (resp. v').

Computing P_2 for RPA to kRR: Since each fake user engaged in RPA attacks randomly selects a value from the encoded data space $\mathcal{D} = \{1, 2, \dots, d\}$, each item $v \in \mathcal{D}$ is selected in each round with probability $\frac{1}{d}$. Therefore, the first and second rounds select the same item v with probability $\frac{1}{d} * \frac{1}{d}$. Furthermore, since v could be any element in \mathcal{D} ,

$$\begin{aligned} P_2 &= \sum_{v \in \mathcal{D}} Pr(v \text{ in 1st round}) * Pr(v \text{ in 2nd round}) \\ &= \sum_{v \in \mathcal{D}} \frac{1}{d} * \frac{1}{d} = d * \left(\frac{1}{d} * \frac{1}{d} \right) = \frac{1}{d}. \end{aligned}$$

Computing P_2 for RIA to kRR: As a fake user is to select a item v from target set $T = \{t_1, t_2, \dots, t_r\}$ and perturbs it to a different value in \mathcal{D} (resp. remains the same value) with probability $q' = q_1$ or q_2 (resp. $p' = p_1$ or p_2), computing P_2 for RIA to kRR should be divided into two cases. The first is that the values reported in the two rounds are the same and belong to T . In this case, the probability that the fake user reports the same value is $\sum_{v \in T} (\frac{1}{r}p' + (1 - \frac{1}{r})q')^2$ where $\frac{1}{r}p'$ (resp. $(1 - \frac{1}{r})q'$) is the probability that the reported value is not perturbed (resp. perturbed). The second is that the values reported in two rounds are the same, and not in T . In this case, the reported value must be perturbed from the true value. Thus, the probability that the fake user reports the same value is $\sum_{v \in \mathcal{D} \setminus T} (q')^2$. Therefore,

$$\begin{aligned} P_2 &= \sum_{v \in T} Pr(v \text{ in 1st round}) * Pr(v \text{ in 2nd round}) \\ &\quad + \sum_{v \in \mathcal{D} \setminus T} Pr(v \text{ in 1st round}) * Pr(v \text{ in 2nd round}) \\ &= \sum_{v \in T} \left(\frac{1}{r}p' + \left(1 - \frac{1}{r}\right)q' \right)^2 + \sum_{v \in \mathcal{D} \setminus T} (q')^2 \\ &= r \left(\frac{1}{r}p' + \left(1 - \frac{1}{r}\right)q' \right)^2 + (d - r)(q')^2. \end{aligned}$$

Computing P_2 for MGA to kRR: Unlike RPA to kRR, each fake user randomly selects a value from target set T instead of encoded data space \mathcal{D} , each item $v \in T$ is selected in each round with a probability $\frac{1}{r}$. Therefore, each fake user selects the same item v in both the first and second rounds with probability $\frac{1}{r} * \frac{1}{r}$. Since v could be any item in T ,

$$\begin{aligned} P_2 &= \sum_{v \in T} Pr(v \text{ in 1st round}) * Pr(v \text{ in 2nd round}) \\ &= \sum_{v \in T} \frac{1}{r} * \frac{1}{r} = r * \left(\frac{1}{r} * \frac{1}{r} \right) = \frac{1}{r}. \end{aligned}$$

B. Percentage Estimation for OUE

For OUE, each genuine user encodes his value v to a d -dimensional binary vector B where the v -th bit is 1, i.e., $\text{Encode}(v) = B$ where $B[v] = 1$ and $B[v'] = 0$ for $v' \neq v$. Each non-zero entry (resp. zero entry) in B is kept unchanged with probability $p' = p_1$ or p_2 (resp. $1 - q'$) and flipped with

probability $1 - p'$ (resp. $q' = q_1$ or q_2). Let B_1 and B_2 be binary vectors reported in the first and second rounds, respectively, $B_1 = B_2$ requires that two corresponding bits in the vectors are the same, and each bit i keeps the same in two rounds when both $B_1[i]$ and $B_2[i]$ are perturbed from $B[i]$ or kept as $B[i]$. For $i = v$, the probability of $B_1[i] = B_2[i]$ is $(p')^2 + (1 - p')^2$ where $(p')^2$ is for $B_1[i] = B_2[i] = B[i]$ and $(1 - p')^2$ is for $B_1[i] = B_2[i] \neq B[i]$. For $i \neq v$, the probability of $B_1[i] = B_2[i]$ is $(q')^2 + (1 - q')^2$. Therefore, the probability P_1 that the data collector collects the same value from a genuine user is

$$\begin{aligned} P_1 &= \sum_{i=v} Pr(B_1[i] = B_2[i]) * \sum_{i \neq v} Pr(B_1[i] = B_2[i]) \\ &= \sum_{i=v} ((p')^2 + (1 - p')^2) * \sum_{i \neq v} ((q')^2 + (1 - q')^2) \\ &= ((p')^2 + (d - 1)(q')^2) ((q')^2 + (1 - q')^2)^{d-1}. \quad (17) \end{aligned}$$

Observe that the dimension of a binary vector could be very large, so P_1 has a very small value and scarce observations of the same value are obtained in two rounds. To solve this problem, P_1 for any τ bits of B_1 and B_2 (denoted by $P_1(\tau)$) should be given. Since the τ bits are randomly selected from d bits, they contain v -th bit with probability $\frac{\tau}{d}$. As such,

$$\begin{aligned} P_1(\tau) &= \frac{\tau}{d} * \sum_{i=v} Pr(B_1[i] = B_2[i]) * \sum_{i \neq v} Pr(B_1[i] = B_2[i]) \\ &= + \left(1 - \frac{\tau}{d}\right) * \sum_{i \neq v} Pr(B_1[i] = B_2[i]) \\ &= \frac{\tau}{d} * ((p')^2 + (d - 1)(q')^2) ((q')^2 + (1 - q')^2)^{\tau-1} \\ &= + \left(1 - \frac{\tau}{d}\right) ((q')^2 + (1 - q')^2)^{\tau}. \end{aligned}$$

Computing P_2 for RPA to OUE: Similar to the computation of P_2 for RPA to kRR, the fake user randomly selects a value from the encoded data space, the difference lie in that the encoded data space is now $\mathcal{D} = \{0, 1\}^d$. Thus, the fake user selects the same binary vector B in the first and second rounds with probability $\frac{1}{2^d} * \frac{1}{2^d}$,

$$\begin{aligned} P_2 &= \sum_{B \in \mathcal{D}} Pr(B \text{ in 1st round}) * Pr(B \text{ in 2nd round}) \\ &= \sum_{B \in \mathcal{D}} \frac{1}{2^d} * \frac{1}{2^d} = 2^d * \left(\frac{1}{2^d} * \frac{1}{2^d}\right) = \frac{1}{2^d}. \end{aligned}$$

Consequently, $P_2(\tau)$ is $2^\tau * (\frac{1}{2^\tau} * \frac{1}{2^\tau}) = \frac{1}{2^\tau}$.

Computing P_2 for RIA to OUE: For RIA to OUE, the fake user first selects a value from target set T and then adopts the same encoding and perturbation procedures used by the genuine users. Let B_1 and B_2 be binary vectors reported in the first and second rounds, respectively, $B_1 = B_2$ requires that two corresponding bits in the vectors are the same, and each bit i keeps the same in two rounds when both $B_1[i]$ and $B_2[i]$ are perturbed from $B[i]$ or kept as $B[i]$. For $i \in T$, the probability of $B_1[i] = B_2[i]$ should be divided into two cases. First, the probability of $B_1[i] = B_2[i] = 1$ is $(\frac{1}{r}p' + (1 - \frac{1}{r})q')^2$. Second, for $B_1[i] = B_2[i] =$

0, the probability is $(\frac{1}{r}(1 - p') + (1 - \frac{1}{r})(1 - q'))^2$. For $i \notin T$, the probability of $B_1[i] = B_2[i]$ is $(q')^2 + (1 - q')^2$. Therefore, the probability P_2 that the data collector collects the same value from a fake user is

$$\begin{aligned} P_2 &= \sum_{i \in T} Pr(B_1[i] = B_2[i]) * \sum_{i \notin T} Pr(B_1[i] = B_2[i]) \\ &= \sum_{i \in T} \left(\frac{1}{r}p' + \left(1 - \frac{1}{r}\right)q' \right)^2 + \left(\frac{1}{r}(1 - p') + \left(1 - \frac{1}{r}\right)(1 - q') \right)^2 * \sum_{i \notin T} ((q')^2 + (1 - q')^2) \\ &= \left[\left(\frac{1}{r}p' + \left(1 - \frac{1}{r}\right)q' \right)^2 + \left(\frac{1}{r}(1 - p') + \left(1 - \frac{1}{r}\right)(1 - q') \right)^2 \right]^r * \left[(q')^2 + (1 - q')^2 \right]^{d-r}. \end{aligned}$$

The computation of P_2 for any chosen τ bits (i.e., $P_2(\tau)$) depends on how many bits (denoted by j) are chosen from T . Therefore, $P_2(\tau)$ is the expectation over j , that is,

$$\begin{aligned} &\sum_{j=0}^r \prod_{i=0}^d Pr(B_1[i] = B_2[i], j) \\ &= \sum_{j=0}^r Pr(j) \prod_{i=0}^d Pr(B_1[i] = B_2[i] | j) \\ &= \sum_{j=0}^r Pr(j) \left[\prod_{i \in T} Pr(B_1[i] = B_2[i] | j) * \prod_{i \notin T} Pr(B_1[i] = B_2[i] | j) \right] \\ &= \sum_{j=0}^r \frac{C_r^j C_{d-r}^{\tau-j}}{C_d^\tau} \left[\prod_{i \in T} Pr(B_1[i] = B_2[i] | j) * \prod_{i \notin T} Pr(B_1[i] = B_2[i] | j) \right] \\ &= \sum_{j=0}^r \frac{C_r^j C_{d-r}^{\tau-j}}{C_d^\tau} \left[\left[\left(\frac{p'}{r} + \frac{r-1}{r}q' \right)^2 + \left(\frac{1}{r}(1 - p') + \frac{r-1}{r}(1 - q') \right)^2 \right]^j * \left[(q')^2 + (1 - q')^2 \right]^{\tau-j} \right] \end{aligned}$$

where $C_{X_1}^{X_2}$ is the number of X_2 -combinations of the set with X_1 elements.

Computing P_2 for MGA to OUE: Recall that each fake user attacks OUE with MGA by sending a binary vector where all target bits and randomly chosen $l = \lfloor p + (d - 1)q - r \rfloor$ bits are 1, and the remaining bits are 0. Let B_1 and B_2 be binary vectors reported in the first and second rounds, respectively, whether or not $B_1 = B_2$ depends on the randomly chosen l bits in two

different rounds. Since these l bits are randomly selected from all non-target bits, each possible l bits are selected with probability of $1/C_{d-r}^l$. Let \mathbb{B}_ℓ be all possible binary vectors where all target bits and additional l bits are 1,

$$\begin{aligned} P_2 &= \sum_{B \in \mathbb{B}_\ell} Pr(B \text{ in 1st round}) * Pr(B \text{ in 2nd round}) \\ &= \sum_{B \in \mathbb{B}_\ell} \frac{1}{C_{d-r}^l} * \frac{1}{C_{d-r}^l} = C_{d-r}^l * \left(\frac{1}{C_{d-r}^l} * \frac{1}{C_{d-r}^l} \right) = \frac{1}{C_{d-r}^l}. \end{aligned}$$

Observe that $P_2(\tau)$ for any chosen τ bits depends on how many bits (denoted by j) are chosen from T . In particular, if j bits are chosen from T , the values in the corresponding entries in B_1 and B_2 are always the same, thus, $P_2(\tau)$ is computed as follows:

$$\begin{aligned} &\sum_{j=0}^r \prod_{i=0}^d Pr(B_1[i] = B_2[i], j) \\ &= \sum_{j=0}^r Pr(j) \prod_{i=0}^d Pr(B_1[i] = B_2[i] | j) \\ &= \sum_{j=0}^r Pr(j) \left[\prod_{i \in T} Pr(B_1[i] = B_2[i] | j) \right. \\ &\quad \left. * \prod_{i \notin T} Pr(B_1[i] = B_2[i] | j) \right] \\ &= \sum_{j=0}^r \frac{C_r^j C_{d-r}^{\tau-j}}{C_d^\tau} \left[1 * \prod_{i \notin T} Pr(B_1[i] = B_2[i] | j) \right] \\ &= \sum_{j=0}^r \frac{C_r^j C_{d-r}^{\tau-j}}{C_d^\tau} \left[\left(\frac{l}{d-r} \right)^2 + \left(\frac{d-r-l}{d-r} \right)^2 \right]^{\tau-j} \end{aligned}$$

where $C_{X_1}^{X_2}$ is the number of X_2 -combinations of the set with X_1 elements, and $(\frac{l}{d-r})^2$ (resp. $(\frac{d-r-l}{d-r})^2$) is the corresponding probability that the bit i is in (resp. not in) the randomly chosen l bits.

C. Percentage Estimation for OLH

We follow Section IV-B3 to set $p^* = \frac{e^\epsilon}{e^\epsilon + d' - 1}$, $q^* = \frac{1}{e^\epsilon + d' - 1}$, $p = p^* = \frac{e^\epsilon}{e^\epsilon + d' - 1}$ and $q = \frac{1}{d'} p^* + (1 - \frac{1}{d'}) q^* = \frac{1}{d'}$. Similarly, p_1^* , q_1^* , p_1 and q_1 (resp. p_2^* , q_2^* , p_2 and q_2) are defined w.r.t. ϵ_1 (resp. ϵ_2). Also, let $p_{1,2}^* = p_1^*$ or p_2^* , $q_{1,2}^* = q_1^*$ or q_2^* . Recall that each genuine user reports his value v with OLH first randomly selects a hash function H from \mathcal{H} to obtained $H(v)$, and then reports the key-value pair $\langle H, H(v)' \rangle$ where $H(v)'$ is perturbed (resp. kept) from $H(v)$ with probability $q_{1,2}^*$ (resp. $p_{1,2}^*$). As such, the genuine user reports the same value when 1) reporting the true value $H(v)$ in both the first and second rounds; 2) reporting the same value $H(v)'$ ($H(v)' \neq H(v)$). The probability P_1 that the data collector collects the same value from a genuine user is

$$P_1 = (p_{1,2}^*)^2 + (d' - 1)(q_{1,2}^*)^2 \quad (18)$$

where $p_{1,2}^*$ (resp. $q_{1,2}^*$) is the probability that the genuine user reports the same value $H(v)$ (resp. $H(v)'$, $H(v)' \neq H(v)$).

Computing P_2 for RPA to OLH: For RPA to OLH, each fake user reports a key-value pair $\langle H, h \rangle$ where h is randomly selected from $\mathcal{D} = \{1, 2, \dots, d'\}$ ($d' = e^\epsilon + 1$) without perturbation. As such, each item $h \in \mathcal{D}$ is selected in each round with probability $\frac{1}{d'}$. Therefore, each fake user selects the same item h in both the first and second rounds with probability $\frac{1}{d'} * \frac{1}{d'}$,

$$\begin{aligned} P_2 &= \sum_{h \in \mathcal{D}} Pr(h \text{ in 1st round}) * Pr(h \text{ in 2nd round}) \\ &= \sum_{h \in \mathcal{D}} \frac{1}{d'} * \frac{1}{d'} = d' * \left(\frac{1}{d'} * \frac{1}{d'} \right) = \frac{1}{d'}. \end{aligned}$$

Computing P_2 for RIA to OLH: Each fake user randomly selects a value from target set T and perturbs it to a value in $\mathcal{D} = \{1, 2, \dots, d'\}$ where $d' = e^\epsilon + 1$, he reports the same value in $\mathcal{D} = \{1, 2, \dots, d'\}$ in two rounds with probability $\frac{1}{d'} * p_{1,2}^* + (1 - \frac{1}{d'}) * q_{1,2}^*$,

$$\begin{aligned} P_2 &= \sum_{h \in \mathcal{D}} Pr(h \text{ in 1st round}) * Pr(h \text{ in 2nd round}) \\ &= \sum_{h \in \mathcal{D}} \left(\frac{1}{d'} * p_{1,2}^* + \left(1 - \frac{1}{d'} \right) * q_{1,2}^* \right)^2 \\ &= d' * \left(\frac{1}{d'} * p_{1,2}^* + \left(1 - \frac{1}{d'} \right) * q_{1,2}^* \right)^2. \end{aligned}$$

Computing P_2 for MGA to OLH: Since MGA to OLH is to seek for a hash function that maps all target values to the same encoded value, P_2 is always be 1. But in practice, there may be no such hash function. Hence, we resort to the hash function that results in the most targets to the same value. Suppose that there are maximally NUM targets hashed to the same value, $P_2 = (\frac{\text{NUM}}{r})^2 + (1 - \frac{\text{NUM}}{r})^2$.

VI. DETECTION OF DATA POISONING ATTACKS

The countermeasures mentioned above assume that the aggregator is aware of the poisoning methods. However, there is a possibility that the aggregator is unaware of the specific poisoning method being used and lacks knowledge about appropriate defense strategies. To overcome this challenge, we have developed a novel method for detecting data poisoning attacks. This method enables the aggregator to identify the attacks and adopt suitable countermeasures to mitigate their negative impact.

The detection method is based on reported values in two rounds and the following observations. First, in the case of the kRR protocol with RIA, the reported values of fake users with RIA as well as genuine users are perturbed, while the values of those under RPA and MGA remain unperturbed. In other words, both fake users and genuine users exhibit similar behavior under the kRR protocol with RIA. Consequently, the frequency distribution of items in two different rounds may appear similar for RIA. Second, RPA and MGA directly manipulate the perturbed values for fake users without utilizing the LDP protocol to generate them. Therefore, for OUE and OLH, the perturbed values of fake users may exhibit statistically abnormal characteristics [25]. For instance, in MGA attacks,

all target items are assigned a value of 1 for each fake user, resulting in a set of binary vectors from fake users where all target items are consistently set to 1. As such, for RIA, the frequency distribution of items in two different rounds may be similar. This can also be observed among genuine users who perturbed their values. Second, RPA and MGA directly craft the perturbed values for fake users, without using the LDP protocol to generate the perturbed values, hence, for OUE and OLH, the perturbed values for the fake users may be statistically abnormal [25]. For example, all target items in MGA attacks are assigned 1 for each fake user, which results in a set of items in binary vectors of the fake users will always be 1.

Algorithm 2 outlines the detection method. Given values Y_1 and Y_2 collected in two rounds and encoded space \mathcal{D} , if the LDP protocol is kRR, it first computes the item frequency distribution $dist_1$ and $dist_2$ for Y_1 and Y_2 (Line 2, Algorithm 2) and measures the similarity of distributions of $dist_1$ and $dist_2$ (i.e., JS_1) using Jensen-Shannon divergence (Line 7). Then, it generates $N + M$ values from \mathcal{D} and adopts the same behavior as genuine users to perturb them to obtain Y'_1 and Y'_2 (Line 3), computes the item frequency distribution $dist'_1$ and $dist'_2$ for Y'_1 and Y'_2 (Line 4), and calculates Jensen-Shannon divergence (i.e., JS_2) of $dist'_1$ and $dist'_2$ (Line 7). Next, it generates $N + M$ values from \mathcal{D} and perturb some of them to obtain Y''_1 and Y''_2 (Line 5), computes the item frequency distribution $dist''_1$ and $dist''_2$ for Y''_1 and Y''_2 (Line 6), and calculates Jensen-Shannon divergence (i.e., JS_3) of $dist''_1$ and $dist''_2$ (Line 7). If JS_1 is closer to JS_2 than JS_3 , it indicates that all values in Y_1 and Y_2 may be perturbed (Lines 8 and 9). This suggests that fake users may adopt the RIA strategy. Otherwise, fake users may adopt either RPA or MGA (Line 10). It is worth noting that differentiating between RPA and MGA is challenging. However, if we cannot differentiate between the two, the countermeasures we adopt may not significantly impact the overall performance.

If the LDP protocol is OUE, we can adapt the frequent itemset mining-based detection method [25] to predict if the attack is RIA or RPA/MGA. Let B be an item set, $z = |B|$, s the number of users whose perturbed binary vectors are 1 for all items in B , and η the maximum false positive rate of detecting fake users, if $s \geq \tau_z$ where τ_z satisfies $\tau_z > (n + m)p'(q')^{z-1}$ and $\frac{(n+m)p'(q')^{z-1}(1-(q')^{z-1})}{[\tau_z - (n+m)p'(q')^{z-1}]^2} \leq \eta$, the attack is predicted as RPA or MGA (Lines 12 to 14). Otherwise, it is predicted as RIA. It further differentiates between RPA and MGA by checking the gap between the maximum and minimum number of 1s in the reported binary vectors. If the gap is larger than half the length of the binary vector, the attack is predicted as RPA. This is because there may be a significant discrepancy in the number of 1s among binary vectors from different fake users employing RPA, who randomly report binary values from the large space of 2^d , where d is the number of items.

For the OLH protocol, fake users employing the MGA strategy aim to hash all targets to the same value. To determine if the attack is MGA, we can represent the hash value a in a tuple (H, a) as a binary vector B , where $B[v] = 1$ if and only if $H(v) = a$. We can apply a similar frequent itemset mining method as before to detect the MGA attack. Specifically, if the false positive rate is at most η , the attack is classified as MGA if

Algorithm 2: Detection of Data Poisoning Attacks.

Input: Y_1, Y_2 , encoded space \mathcal{D}
Output: detected poisoning attacks

- 1: **if** LDP protocol is kRR **then**
- 2: $dist_1 \leftarrow \text{ITEMFREQUENCY}(Y_1)$,
 $dist_2 \leftarrow \text{ITEMFREQUENCY}(Y_2)$
- 3: generate $N+M$ values from \mathcal{D} and perturb them to Y'_1 , Y'_2
- 4: $dist'_1 \leftarrow \text{ITEMFREQUENCY}(Y'_1)$,
 $dist'_2 \leftarrow \text{ITEMFREQUENCY}(Y'_2)$
- 5: generate $N+M$ values from \mathcal{D} , perturb some of them and report Y''_1, Y''_2
- 6: $dist''_1 \leftarrow \text{ITEMFREQUENCY}(Y''_1)$,
 $dist''_2 \leftarrow \text{ITEMFREQUENCY}(Y''_2)$
- 7: $JS_1 = \text{JSD}(dist_1, dist_2)$, $JS_2 = \text{JSD}(dist'_1, dist'_2)$,
 $JS_3 = \text{JSD}(dist''_1, dist''_2)$
- 8: **if** $|JS_1 - JS_2| \leq |JS_1 - JS_3|$ **then**
- 9: **return** RIA
- 10: **else return** RPA or MGA
- 11: **Else if** LDP protocol is OUE
- 12: **For** $B \in \text{ALLITEMS}(Y_1)$
- 13: **if** $s \geq \tau_z$ where $\tau_z > (n + m)p'(q')^{z-1}$ and
 $\frac{(n+m)p'(q')^{z-1}(1-(q')^{z-1})}{[\tau_z - (n+m)p'(q')^{z-1}]^2} \leq \eta$ **then**
- 14: **return** RPA or MGA
- 15: **return** RIA
- 16: **Else if** DP protocol is OLH
- 17: **For** $B \in \text{ALLITEMS}(\text{VECTORIZATION}(Y_1))$
- 18: **if** $I(q'; \tau_z, n + m - \tau_z + 1) \leq \eta$ **then return** MGA
- 19: **return** RIA or RPA

a binary vector B satisfies $I(q'; \tau_z, n + m - \tau_z + 1) \leq \eta$ [25], where I denotes the regularized incomplete beta function [64] (Lines 16 to 18). Otherwise, the attack is either RIA or RPA. Although we cannot differentiate between RPA and RIA, the choice of countermeasures for RPA and RIA may not significantly impact the overall performance.

VII. EXPERIMENTAL EVALUATION

In this section, we investigate the performance of LDPGuard and report the key findings. LDPGuard is implemented with Python 3.10. All experiments are conducted on macOS Monterey 12.2.1 with Quad-Core Intel Core i5 4-Core CPU (2 GHz) and 16 GB of main memory.

A. Experimental Setup

Datasets: LDPGuard is evaluated on three datasets consisting of two real-world datasets (i.e., IPUMS [62] and Fire [63]) and one synthetic dataset that follows the Zipf distribution.

Baselines: We compare LDPGuard with two related countermeasures [25] against data poisoning attacks: (1) *Normalization* (denoted by **NORM**) that normalizes the estimated frequency of the items so that each estimated item frequency is non-negative and the overall summation of item frequencies is 1; (2) *Fake users detection* (denoted by **DETECT**) that detects possible fake

TABLE II
PARAMETER SETTINGS FOR EXPERIMENTS

Parameter	Symbol	Default Value
number of targets	r	10
percentage of fake users	β	0.05
privacy budget	ϵ	1.0
	ϵ_1	$\frac{\epsilon}{2}$
	ϵ_2	$\frac{\epsilon}{2}$

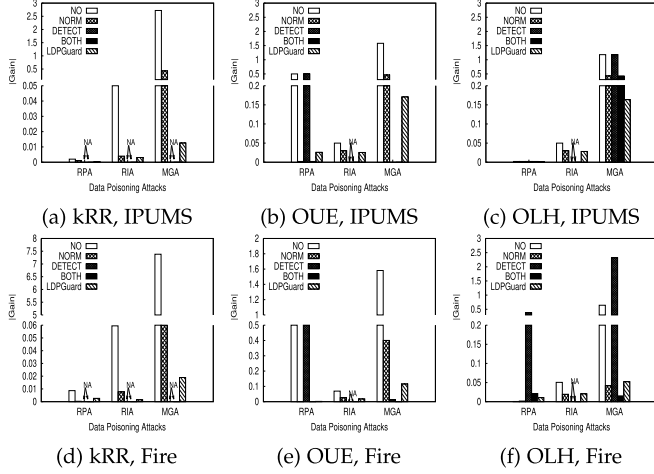


Fig. 2. Absolute Gain $|Gain|$ on IPUMS and Fire.

users based on their reported values to alleviate their impact on the final results. We also compare it with the results where no countermeasure is adopted (denoted by **NO**), and both DETECT and NORM are applied (denoted by **BOTH**). Noted that the resulting frequency of LDPGuard is also normalized to avoid negative results.

Performance Metrics: We use **absolute gain**, i.e., the absolute value of the frequency gain (denoted by $|Gain|$), to compare LDPGuard with baselines.¹ We use **NA** to denote the case that a countermeasure is not applicable for a data poisoning attack. We compare LDPGuard with DETECT in terms of the estimated percentage (i.e., $\hat{\beta}$) of fake users.

Parameter Settings: Unless otherwise stated, we use the default parameter values in Table II. In addition, we follow existing work [13], [25] to set $d' = e^\epsilon + 1$ for OLH.

B. Experimental Results

Exp 1. Performances of Frequency Estimation: We first compare LDPGuard with baselines on three datasets and plot the results of frequency estimation on IPUMS and Fire in Fig. 2. As depicted in Fig. 2(a) and (d), LDPGuard always outperforms baselines for kRR regardless of data poisoning attacks. Compared to NO, NORM reduces $|Gain|$ to a relatively smaller value, but the impact is limited as NORM does not reduce

¹One may argue that he can follow the existing work [25] to use frequency gain $Gain$ instead of its absolute value, but he may ignore the fact that the frequency of targets is largely distorted if $Gain$ is a negative number with a large absolute value, which results in unfair comparison for countermeasures.

²Due to space limitations, results on Zipf are given in [1].

TABLE III
PERCENTAGE ESTIMATION ON IPUMS&Fire, $\beta = 0.05$

Protocols	Attacks	$\hat{\beta}$ (IPUMS)		$\hat{\beta}$ (Zipf)	
		DETECT	LDPGuard	DETECT	LDPGuard
KRR	RPA	—	0.117	—	0.593
	RIA	—	0.129	—	0.382
	MGA	—	0.050	—	0.049
OUE	RPA	7.0E-05	0.060	7.0E-05	0.042
	RIA	—	0.232	—	0.292
	MGA	0.050	0.044	0.050	0.041
OLH	RPA	0.029	0.013	0.025	0.081
	RIA	—	0.024	—	0.045
	MGA	0.346	0.048	0.267	0.040

$|Gain|$ directly but normalizes each estimated item frequency to a non-negative value. For kRR, both DETECT and BOTH cannot work and their results are denoted by NA. For OUE (Fig. 2(b) and (e)) and OLH (Fig. 2(c) and (f)), NORM is ineffective in the scenario where MGA attacks are deployed, and the reason is discussed above. DETECT and BOTH are effective in this scenario, as they can detect abnormal statistical patterns in the reported values generated by MGA. However, DETECT is ineffective for RPA because each fake user in RPA randomly selects perturbed values in the encoded space, and thus the perturbed values do not show meaningful statistical patterns. In addition, both DETECT and BOTH are not applicable for RIA since each fake user under RIA perturbs the value sampled from the target item set before sending it to a data collector, which results in the incapability of distinguishing fake users from genuine users. In contrast, LDPGuard can obtain less $|Gain|$ than NORM for MGA. Compared to DETECT, it has better performance for RIA. Moreover, LDPGuard is applicable regardless of LDP protocols and data poisoning attacks.

In summary, LDPGuard can effectively defend against data poisoning attacks and outperforms baselines in terms of absolute gain.

Exp 2. Performances of Percentage Estimation: We further compare LDPGuard with DETECT on IPUMS and Fire,³ and report the estimated percentages of fake users in Table III. In this experiment, we set the true percentage of fake users to 0.05 and use “—” to denote the case where the countermeasure is not applicable. Observe that DETECT can accurately estimate the percentage of MGA to OUE because each reported value of OUE is a binary vector and the target bits in the vector tends to expose abnormal statistical patterns. We also observe that it has limitations in estimating the percentage in the scenarios where there is no deterministic and meaningful statistical patterns. For example, for RPA to OUE, each fake user randomly selects a value from a very large encoding space \mathcal{D} (i.e., $\mathcal{D} = \{0, 1\}^d$) and reports it without perturbation. This will result in each vector reported by fake users being with about $d/2$ non-zero bits, and further lead to the difficulty of distinguishing a fake user from a genuine user who reports the value with nearly $1/2 + (d - 1)/(e + 1)$ non-zero bits when $\epsilon = 1$. Furthermore, DETECT is not applicable for KRR as well as OUE (and OLH) to RIA.

³Due to space limitations, results on Zipf are given in [1].

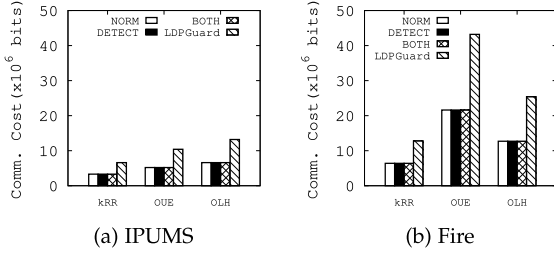


Fig. 3. Communication cost.

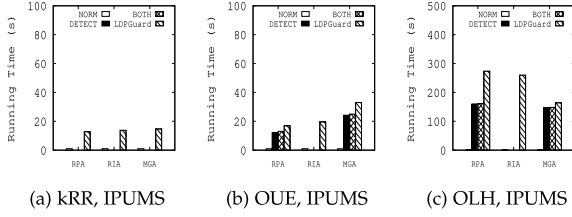


Fig. 4. Computation cost.

TABLE IV

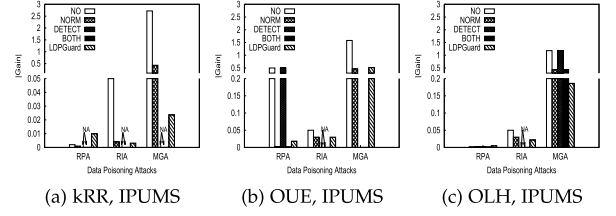
 PERCENTAGE ESTIMATION IN UNKNOWN ATTACK SCENARIOS, $\beta = 0.05$

Protocols	Attacks	$\tilde{\beta}$ (IPUMS)	
		DETECT	LDPGuard
KRR	RPA	—	0.001
	RIA	—	0.002
	MGA	—	0.048
OUE	RPA	7.0E-05	0.009
	RIA	—	0.001
	MGA	0.050	0.014
OLH	RPA	0.029	0.001
	RIA	—	0.053
	MGA	0.346	0.047

In contrast, LDPGuard is applicable for all LDP protocols and data poisoning attacks, although the estimated percentages under RPA and RIA where randomnesses are here and there are not very precise. In general, the estimated percentages by LDPGuard are close to the true percentages (e.g., 0.05) of fake users. Hence, LDPGuard is able to accurately estimate the percentage of fake users.

Exp 3. Computational and Communication Efficiency: We also conducted a comparison of computational and communication efficiency. Fig. 3 illustrates the communication cost. Compared to competitors, LDPGuard incurs higher communication costs. This is because LDPGuard requires two rounds of data collection. We also compared their running times in Fig. 4. If the running time is 0 s, it means that the corresponding method does not support this LDP protocol. We observed that LDPGuard requires relatively more time because it simulates attacker behavior to counter the attack, but the time taken by LDPGuard remains limited to minutes.

Exp 4. Performance in Unknown Attack Scenarios: We further compared the performance in unknown attack scenarios. From Table IV and Fig. 5, we can draw the same conclusions as in Exp 1 and Exp 2. Specifically, Table IV indicates that even though we assumed DETECT to have knowledge of the attack types


 Fig. 5. Absolute gain $|Gain|$.

in this experiment, LDPGuard still demonstrates significant effectiveness in estimating the percentage of fake users, albeit with slightly lower accuracy, while DETECT performs well only against MGA for OUE and RPA for OLH. Additionally, Fig. 5 shows that LDPGuard is capable of mitigating the impact of attacks even in cases where the attack type is unknown in advance.

VIII. CONCLUSION

In this article, we focus on the problem of defending against data poisoning attacks to local differential privacy protocols. We present a novel framework called LDPGuard to this end. In particular, LDPGuard first adopts the two-round collection strategy to accurately estimate the percentage of fake users. Then it utilizes the estimated percentage to provide adversarial schemes to defend against various data poisoning attacks. LDPGuard supports state-of-the-art LDP protocols (kRR, OUE, and OLH) for frequency estimation and can effectively defend against existing data poisoning attacks. Experimental study on real-world and synthetic datasets demonstrates the superiority of LDPGuard compared to existing techniques. As part of our future work, we will explore data poisoning attacks and countermeasures to other data types, such as graph data.

REFERENCES

- [1] K. Huang, et al., "LDPGuard: Defenses against data poisoning attacks to local differential privacy protocols," 2023, pp. 1–17. [Online]. Available: <https://github.com/TechReport2023/LDPGuard/blob/main/LDPGuard-TR.pdf>
- [2] B. Avent, A. Korolova, D. Zeber, T. Hovden, and B. Livshits, "BLENDER: Enabling local search with a hybrid differential privacy model," in *Proc. USENIX Secur. Symp.*, 2017, pp. 747–764.
- [3] R. Bassily, K. Nissim, U. Stemmer, and A. G. Thakurta, "Practical locally private heavy hitters," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2017, pp. 2285–2293.
- [4] R. Bassily and A. Smith, "Local, private, efficient protocols for succinct histograms," in *Proc. Conf. Symp. Theory Comput.*, 2015, pp. 127–135.
- [5] G. Cormode, T. Kulkarni, and D. Srivastava, "Marginal release under local differential privacy," in *Proc. Int. Conf. Manage. Data*, 2018, pp. 131–146.
- [6] J. C. Duchi, M. I. Jordan, and M. J. Wainwright, "Local privacy and statistical minimax rates," in *Proc. IEEE Annu. Symp. Foundations Comput. Sci.*, 2013, pp. 429–438.
- [7] Ú. Erlingsson, V. Pihur, and A. Korolova, "RAPPOR: Randomized aggregatable privacy-preserving ordinal response," in *Proc. ACM Asia Conf. Comput. Commun. Secur.*, 2014, pp. 1054–1067.
- [8] J. Jia and N. Z. Gong, "Calibrate: Frequency estimation and heavy hitter identification with local differential privacy via incorporating prior knowledge," in *Proc. IEEE Conf. Comput. Commun.*, 2019, pp. 2008–2016.
- [9] P. Kairouz, K. Bonawitz, and D. Ramage, "Discrete distribution estimation under local privacy," in *Proc. Int. Conf. Mach. Learn.*, 2016, pp. 2436–2444.

- [10] P. Kairouz, S. Oh, and P. Viswanath, "Extremal mechanisms for local differential privacy," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2014, pp. 2879–2887.
- [11] Z. Qin, Y. Yang, T. Yu, I. Khalil, X. Xiao, and K. Ren, "Heavy hitter estimation over set-valued data with local differential privacy," in *Proc. ACM Asia Conf. Comput. Commun. Secur.*, 2016, pp. 192–203.
- [12] X. Ren et al., "LoPub: High-dimensional crowdsourced data publication with local differential privacy," *IEEE Trans. Inf. Forensics Secur.*, vol. 13, no. 9, pp. 2151–2166, Sep. 2018.
- [13] T. Wang, J. Blocki, N. Li, and S. Jha, "Locally differentially private protocols for frequency estimation," in *Proc. USENIX Secur. Conf.*, 2017, pp. 729–745.
- [14] T. Wang et al., "Answering multi-dimensional analytical queries under local differential privacy," in *Proc. Int. Conf. Manage. Data*, 2019, pp. 159–176.
- [15] T. Wang, N. Li, and S. Jha, "Locally differentially private frequent itemset mining," in *Proc. IEEE Eur. Symp. Secur. Privacy*, 2018, pp. 127–143.
- [16] T. Wang, N. Li, and S. Jha, "Locally differentially private heavy hitter identification," *IEEE Trans. Dependable Secure Comput.*, vol. 18, no. 2, pp. 982–993, Mar./Apr. 2021.
- [17] T. Wang, M. Lopusha-Zwakenberg, Z. Li, B. Skoric, and N. Li, "Locally differentially private frequency estimation with consistency," in *Proc. Annu. Netw. Distrib. Syst. Secur. Symp.*, 2020, pp. 1–16.
- [18] Q. Ye et al., "PrivKVM*: Revisiting key-value statistics estimation with local differential privacy," *IEEE Trans. Dependable Secure Comput.*, vol. 20, no. 1, pp. 17–35, Jan./Feb. 2022.
- [19] Z. Zhang, T. Wang, N. Li, S. He, and J. Chen, "CALM: Consistent adaptive local marginal for marginal release under local differential privacy," in *Proc. ACM Asia Conf. Comput. Commun. Secur.*, 2018, pp. 212–229.
- [20] Apple Differential Privacy Team, "Learning with privacy at scale," 2017. [Online]. Available: <https://machinelearning.apple.com/research/learning-with-privacy-at-scale>
- [21] Q. Ye, H. Hu, X. Meng, and H. Zheng, "PrivKV: Key-value data collection with local differential privacy," in *Proc. IEEE Eur. Symp. Secur. Privacy*, 2019, pp. 317–331.
- [22] Q. Ye, H. Hu, M. H. Au, X. Meng, and X. Xiao, "LF-GDPR: A framework for estimating graph metrics with local differential privacy," *IEEE Trans. Knowl. Data Eng.*, vol. 34, no. 10, pp. 4905–4920, Oct. 2020.
- [23] B. Ding, J. Kulkarni, and S. Yekhanin, "Collecting telemetry data privately," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2017, pp. 3574–3583.
- [24] T. T. Nguyen et al., "Collecting and analyzing data from smart device users with local differential privacy," 2016, *arXiv:1606.05053*.
- [25] X. Cao, J. Jia, and N. Z. Gong, "Data poisoning attacks to local differential privacy protocols," in *Proc. USENIX Secur. Conf.*, 2021, pp. 947–964.
- [26] L. Sweeney, "k-anonymity: A model for protecting privacy," *Int. J. Uncertainty Fuzziness Knowl.-Based Syst.*, vol. 10, pp. 557–570, 2002.
- [27] M. Ashwin et al., "L-diversity: Privacy beyond k-anonymity," in *Proc. Int. Conf. Data Eng.*, 2006, pp. 24–24.
- [28] N. Li, T. Li, and V. Suresh, "t-closeness: Privacy beyond k-anonymity and l-diversity," in *Proc. Int. Conf. Data Eng.*, 2007, pp. 106–115.
- [29] C. Zhao, L. Zou, and F. Li, "Privacy preserving subgraph matching on large graphs in cloud," in *Proc. Int. Conf. Manage. Data*, 2016, pp. 199–213.
- [30] C. Dwork, F. McSherry, K. Nissim, and A. Smith, "Calibrating noise to sensitivity in private data analysis," in *Proc. Theory Cryptogr. Conf.*, 2006, pp. 265–284.
- [31] Q. Ye, H. Hu, N. Li, X. Meng, H. Zheng, and H. Yan, "Beyond value perturbation: Local differential privacy in the temporal setting," in *Proc. IEEE Conf. Comput. Commun.*, 2021, pp. 1–10.
- [32] K. Huang, H. Hu, S. Zhou, J. Guan, Q. Ye, and X. Zhou, "Privacy and efficiency guaranteed social subgraph matching," *Vldb J.*, vol. 31, pp. 581–602, 2021.
- [33] C. Dwork and A. Roth, "The algorithmic foundations of differential privacy," *Foundations Trends Theor. Comput. Sci.*, vol. 9, pp. 211–407, 2014.
- [34] N. Li, M. Lyu, D. Su, and W. Yang, "Differential privacy: From theory to practice," in *Synthesis Lectures Information Security, Privacy, and Trust*, Berlin, Germany: Springer, 2016.
- [35] S. L. Warner, "Randomized response: A survey technique for eliminating evasive answer bias," *J. Amer. Stat. Assoc.*, vol. 60, pp. 63–66, 1965.
- [36] S. P. Kasiviswanathan et al., "What can we learn privately?," *SIAM J. Comput.*, vol. 40, pp. 793–826, 2011.
- [37] J. Hsu, S. Khanna, and A. Roth, "Distributed private heavy hitters," in *Automata, Languages, Programming*, Berlin, Germany: Springer, 2012, pp. 461–472.
- [38] G. Fanti, V. Pihur, and Ú. Erlingsson, "Building a RAPPOR with the Unknown: Privacy-preserving learning of associations and data dictionaries," in *Proc. Annu. Privacy Enhancing Technol. Symp.*, 2016, pp. 1–17.
- [39] A. Cheu, A. Smith, and J. Ullman, "Manipulation attacks in local differential privacy," in *Proc. IEEE Eur. Symp. Secur. Privacy*, 2021, pp. 883–900.
- [40] S. Alfeld, X. Zhu, and P. Barford, "Data poisoning attacks against autoregressive models," in *Proc. AAAI Conf. Artif. Intell.*, 2016, pp. 1452–1458.
- [41] B. Biggio, B. Nelson, and P. Laskov, "Poisoning attacks against support vector machines," in *Proc. Int. Conf. Mach. Learn.*, 2012, pp. 1467–1474.
- [42] M. Fang, X. Cao, J. Jia, and N. Gong, "Local model poisoning attacks to byzantine robust federated learning," in *Proc. USENIX Secur. Conf.*, 2020, pp. 1605–1622.
- [43] M. Jagielski, A. Oprea, B. Biggio, C. Liu, C. Nita-Rotaru, and B. Li, "Manipulating machine learning: Poisoning attacks and countermeasures for regression learning," in *Proc. IEEE Eur. Symp. Secur. Privacy*, 2018, pp. 19–35.
- [44] T. Gu, K. Liu, B. Dolan-Gavitt, and S. Garg, "BadNets: Evaluating backdooring attacks on deep neural networks," *IEEE Access*, vol. 7, pp. 47230–47244, 2019.
- [45] A. Shafahi et al., "Poison frogs! Targeted clean-label poisoning attacks on neural networks," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2018, pp. 6106–6116.
- [46] Y. Liu et al., "Trojaning attack on neural networks," in *Proc. Annu. Netw. Distrib. Syst. Secur. Symp.*, 2018, pp. 1–15.
- [47] M. Fang, N. Z. Gong, and J. Liu, "Influence function based data poisoning attacks to top-n recommender systems," in *Proc. World Wide Web Conf.*, 2020, pp. 3019–3025.
- [48] M. Fang, G. Yang, N. Z. Gong, and J. Liu, "Poisoning attacks to graph-based recommender systems," in *Proc. Annu. Comput. Secur. Appl. Conf.*, 2018, pp. 381–392.
- [49] B. Li, Y. Wang, A. Singh, and Y. Vorobeychik, "Data poisoning attacks on factorization-based collaborative filtering," in *Proc. Int. Conf. Neural Inf. Process. Syst.*, 2016, pp. 1893–1901.
- [50] G. Yang, N. Z. Gong, and Y. Cai, "Fake co-visitation injection attacks to recommender systems," in *Proc. Annu. Netw. Distrib. Syst. Secur. Symp.*, 2017, pp. 30–40.
- [51] M. Mozaffari-Kermani, S. Sur-Kolay, A. Raghunathan, and N. K. Jha, "Systematic poisoning attacks on and defenses for machine learning in healthcare," *IEEE J. Biomed. Health Inform.*, vol. 19, no. 6, pp. 1893–1905, Nov. 2015.
- [52] Q. Cao, X. Yang, J. Yu, and C. Palow, "Uncovering large groups of active malicious accounts in online social networks," in *Proc. ACM Asia Conf. Comput. Commun. Secur.*, 2014, pp. 477–488.
- [53] G. Danezis and P. Mittal, "Sybilinifer: Detecting sybil nodes using social networks," in *Proc. Annu. Netw. Distrib. Syst. Secur. Symp.*, 2009, pp. 1–15.
- [54] N. Z. Gong, M. Frank, and P. Mittal, "SybilBelief: A semi-supervised learning approach for structure-based sybil detection," *IEEE Trans. Inf. Forensics Secur.*, vol. 9, no. 6, pp. 976–987, Jun. 2014.
- [55] G. Stringhini, C. Kruegel, and G. Vigna, "Detecting spammers on social networks," in *Proc. 26th Annu. Comput. Secur. Appl. Conf.*, 2010, pp. 1–9.
- [56] B. Wang, J. Jia, and N. Z. Gong, "Graph-based security and privacy analytics via collective classification with joint weight learning and propagation," in *Proc. Annu. Netw. Distrib. Syst. Secur. Symp.*, 2019, pp. 1–15.
- [57] G. Wang, T. Konolige, C. Wilson, X. Wang, H. Zheng, and B. Y. Zhao, "You are how you click: Clickstream analysis for sybil detection," in *Proc. USENIX Secur. Conf.*, 2013, pp. 241–256.
- [58] H. Yu, H. Yu, M. Kaminsky, P. B. Gibbons, and A. Flaxman, "SybilGuard: Defending against sybil attacks via social networks," in *Proc. Conf. Appl. Technol. Architectures Protoc. Comput. Commun.*, 2006, pp. 267–278.
- [59] D. Yuan et al., "Detecting fake accounts in online social networks at the time of registrations," in *Proc. ACM Asia Conf. Comput. Commun. Secur.*, 2019, pp. 1423–1438.
- [60] Y. Collet, "xxHash: Extremely fast hash algorithm," 2016. [Online]. Available: <https://github.com/Cyan4973/xxHash>
- [61] K. Thomas, D. McCoy, C. Grier, A. Kolecz, and V. Paxson, "Trafficking fraudulent accounts: The role of the underground market in twitter spam and abuse," in *Proc. USENIX Secur. Conf.*, 2013, pp. 195–210.
- [62] R. Steven et al., "IPUMS USA: Version 9.0 [dataset]," Minneapolis, MN, USA: IPUMS, 2019. [Online]. Available: <https://doi.org/10.18128/D010.V9.0>
- [63] "San francisco fire department calls for service," 2019. [Online]. Available: <http://bit.ly/336sddl>
- [64] M. Abramowitz and I. A. Stegun, "Handbook of mathematical functions with formulas, graphs, and mathematical tables," 10th ed. New York, NY, USA: Academic, 1972.

- [65] H. Liu, J. Jia, and N. Z. Gong, "PoisonedEncoder: Poisoning the unlabeled pre-training data in contrastive," in *Proc. USENIX Secur. Symp.*, 2022, pp. 3629–3645.
- [66] X. Ren et al., "LDP-IDS: Local differential privacy for infinite data streams," in *Proc. Int. Conf. Manage. Data*, 2022, pp. 1064–1077.
- [67] H. Wang et al., "L-SRR: Local differential privacy for location-based services with staircase randomized response," in *Proc. ACM Asia Conf. Comput. Commun. Secur.*, 2022, pp. 2809–2823.
- [68] M. Zhou et al., "Locally differentially private sparse vector aggregation," in *Proc. IEEE Eur. Symp. Secur. Privacy*, 2022, pp. 422–439.
- [69] Y. Zhang et al., "Trajectory data collection with local differential privacy," 2023, *arXiv:2307.09339*.
- [70] Y. Du et al., "LDPTrace: Locally differentially private trajectory synthesis," 2023, *arXiv:2302.06180*.
- [71] X. Li et al., "Local differentially private heavy hitter detection in data streams with bounded memory," 2023, *arXiv:2311.16062*.
- [72] X. Li et al., "Fine-grained poisoning attack to local differential privacy protocols for mean and variance estimation," in *Proc. USENIX Secur. Symp.*, 2023, pp. 1739–1756.
- [73] Y. Wu et al., "Poisoning attacks to local differential privacy protocols for key-value data," in *Proc. USENIX Secur. Conf.*, 2022, pp. 519–536.
- [74] X. Cao et al., "FedRecover: Recovering from poisoning attacks in federated learning using historical information," in *Proc. IEEE Eur. Symp. Secur. Privacy*, 2023, pp. 1366–1383.



Kai Huang received the BEng degree in software engineering from East China Normal University, in 2014, and the PhD degree from the School of Computer Science, Fudan University, in 2020. He is an Assistant Professor with the School of Computer Science and Engineering, Faculty of Innovation Engineering, Macau University of Science and Technology. He was a postdoc with the Department of Computer Science and Engineering, HKUST, under the supervision of Prof. Xiaofang Zhou. His research interests include graph database and privacy-aware data management.



Gaoya Ouyang received the master's degree from the Department of Mathematics and Information Technology, Education University of Hong Kong, Hong Kong, in 2018. She is currently a research assistant with HKUST, under the supervision of Prof. Xiaofang Zhou. Her research interests include local differential privacy and road network.



Qingqing Ye received the PhD degree in computer science from the Renmin University of China, in 2020. She is an Assistant Professor with the Department of Electronic and Information Engineering, Hong Kong Polytechnic University. She has received several prestigious awards, including Highly Cited Paper Award, China National Scholarship, Outstanding Doctoral Dissertation Award, and IEEE S&P Travel Award. Her research interests include data privacy and security, and adversarial machine learning.



Haibo Hu (Senior Member, IEEE) is a professor with the Department of Electronic and Information Engineering, Hong Kong Polytechnic University and the programme leader of BSc (Hons) in information security. His research interests include cybersecurity, data privacy, Internet of Things, and adversarial machine learning. He has published more than 110 research papers in refereed journals, international conferences, and book chapters. As principal investigator, he has received more than 20 million HK dollars of external research grants from Hong Kong and mainland China.

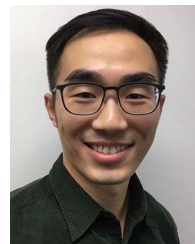
He is an associate editor of *ACM Transactions on Privacy and Security* (TOPS). He is the recipient of a number of titles and awards, including IEEE MDM 2019 Best Paper Award, WAIM Distinguished Young Lecturer, ICDE 2020 Outstanding Reviewer, VLDB 2018 Distinguished Reviewer, ACM-HK Best PhD Paper, Microsoft Imagine Cup, and GS1 Internet of Things Award. He is a CCF, and a certified Cisco CCNA Security Trainer.



Bolong Zheng received bachelor's and master's degrees in computer science from HUST, in 2011 and 2013, respectively, and the PhD degree from the University of Queensland, in 2017. He is an associate professor with the Huazhong University of Science and Technology (HUST). His research interests include spatial-temporal data management. He has published more than 40 papers in prestigious journals and conferences in data management field such as SIGMOD, ICDE, VLDB, ACM Transactions and IEEE Transactions.



Xi Zhao received the master's degree in computer science from the Huazhong University of Science and Technology in 2021. He is currently working toward the PhD degree with the Department of CSE, HKUST, supervised by Prof. Xiaofang Zhou. His research interests include the approximate query processing in high dimensional spaces, exact trajectory similarity search, and exact textual similarity search.



Ruiyuan Zhang received the PhD degree from the University of Queensland, Australia, in 2021. He is currently postdoc fellow with the Hong Kong University of Science and Technology. His research interests include machine learning, and artificial intelligence and their applications in different domain sectors, including energy, logistics, telecommunication, finance, etc.



Xiaofang Zhou (Fellow, IEEE) received the bachelor's and master's degrees in computer science from Nanjing University, in 1984 and 1987 respectively, and the PhD degree in computer science from the University of Queensland, in 1994. He is Otto Poon professor in Engineering and chair professor of Computer Science and Engineering with the Hong Kong University of Science and Technology. His research is focused on finding effective and efficient solutions for managing, integrating, and analysing very large amount of complex data for business, scientific and personal applications. His research interests include spatial and multimedia databases, high performance query processing, web information systems, data mining, data quality management, and machine learning.