# Improving the Utility of Locally Differentially Private Protocols for Longitudinal and Multidimensional Frequency Estimates

Héber H. Arcolezi, Jean-François Couchot
*Femto-ST Institute*
*Univ. Bourgogne Franche-Comté, CNRS*
Belfort, France
heber.hwang_arcolezi@univ-fcomte.fr,
jean-francois.couchot@univ-fcomte.fr

Bechara Al Bouna
*TICKET Lab.*
*Antonine University*
Hadat-Baabda, Lebanon
bechara.albouna@ua.edu.lb

Xiaokui Xiao
*School of Computing*
*National University of Singapore*
Singapore, Singapore
xkxiao@nus.edu.sg

*Abstract*—**This paper investigates the problem of collecting multidimensional data throughout time (i.e., longitudinal studies) for the fundamental task of frequency estimation under local differential privacy (LDP). Contrary to frequency estimation of a single attribute (the majority of the works), the multidimensional aspect imposes to pay particular attention to the privacy budget. Besides, when collecting user statistics longitudinally, privacy progressively degrades. Indeed, both "multiple" settings combined (i.e., many attributes and several collections throughout time) imposes several challenges, in which this paper proposes the first solution for frequency estimates under LDP. To tackle these issues, we extend the analysis of three state-of-the-art LDP protocols (Generalized Randomized Response – GRR, Optimized Unary Encoding – OUE, and Symmetric Unary Encoding – SUE) for both longitudinal and multidimensional data collections. While the known literature uses OUE and SUE for two rounds of sanitization (a.k.a. memoization), i.e., L-OUE and L-SUE, respectively, we analytically and experimentally show that starting with OUE and then with SUE provides higher data utility (i.e., L-OSUE). Also, for attributes with small domain sizes, we propose longitudinal GRR (L-GRR), which provides higher utility than the other protocols based on unary encoding. Lastly, we also propose a new solution named Adaptive LDP for LOngitudinal and Multidimensional FREquency Estimates (ALLOMFREE), which randomly samples a single attribute to send with the whole privacy budget and adaptively selects the optimal protocol, i.e., either L-GRR or L-OSUE. As shown in the results, ALLOMFREE consistently and considerably outperforms the state-of-the-art L-SUE and L-OUE protocols in the quality of the frequency estimations.**

*Index Terms*—**Local differential privacy, Discrete distribution estimation, Frequency estimation, Multidimensional data, Longitudinal studies.**

## I. Introduction

### A. Background

In recent years, differential privacy (DP) [1], [2] has been increasingly accepted as the current standard for data privacy [3]–[6]. With the centralized model of DP, a trusted curator has access to compute on the entire raw data of users (e.g., the Census Bureau [7], [8]). By 'trusted', we mean that curators do not misuse or leak private information from individuals. However, this assumption does not always hold in real life, e.g., data breaches are all too common [9].

To preserve privacy at the user-side, an alternative approach, namely, local differential privacy (LDP), was initially formalized in [10]. With LDP, rather than trusting in a data curator to have the raw data and sanitize it to output queries, each user applies a DP mechanism to their data before transmitting it to the data collector server. The local DP model allows collecting data in unprecedented ways and, therefore, it has led to several adoptions by industry (e.g., Google Chrome browser [11], Microsoft windows 10 operation system [12], Apple iOS and macOS [13]).

### B. Motivation and problem statement

When collecting data in practice, one is often interested in multiple attributes of a population, i.e., *multidimensional data*. For instance, in crowd-sourcing applications, the server may collect both demographic information (e.g., gender, nationality) and user habits in order to develop personalized solutions for specific groups. In addition, one generally aims to collect data from the same users throughout time (i.e., *longitudinal studies*), which is essential in many situations [11], [12]. For example, the fact that two medical acts identified at different times have been performed on the same patient or two different patients means treatment in the first case or two isolated acts in the second.

So, in this paper, we focus on the problem of private frequency (or histogram) estimation of multiple attributes throughout time with LDP. Frequency estimation is a primary objective of LDP, in which the data collector (a.k.a. the aggregator) decodes all the privatized data of the users and can then estimate the number of users for each possible value. More formally, we assume there are $d$ attributes $A = \{A_1, A_2, ..., A_d\}$, where each attribute $A_j$ with a discrete domain $\mathcal{D}_j$ has a specific number of values $k_j = |A_j|$. Each user $u_i$ for $i \in \{1, 2, ..., n\}$ has a tuple $\mathbf{v}^{(i)} = (v_1^{(i)}, v_2^{(i)}, ..., v_d^{(i)})$, where $v_j^{(i)}$ represents the value of attribute $A_j$ in record $\mathbf{v}^{(i)}$. Thus, for each attribute $A_j$ at time $t \in [1, \tau]$, the aggregator's goal

is to estimate a $k_j$-bins histogram, including the frequency of all values in $\mathcal{D}_j$.

Indeed, in both longitudinal and multidimensional settings, one needs to consider the allocation of the privacy budget, which can grow extremely quickly due to the composition theorem [3]. However, on the one hand, most frequency estimation academic literature [14]–[20] focuses on a single data collection (i.e., non-longitudinal studies). On the other hand, the studies for collecting multidimensional data with LDP mainly focused on other complex tasks (e.g., analytical/range queries [21]–[24], estimating marginals [25]–[29]) and numerical data only (e.g., [30]–[33]).

### C. Summary of contributions

In this paper, we extend the analysis of three state-of-the-art LDP protocols, namely, generalized randomized response (GRR) [16], optimized unary encoding (OUE) [14], and symmetric unary encoding (SUE) [11] for both longitudinal and multidimensional frequency estimates. On the one hand, for all three protocols, we theoretically prove that randomly sampling a single attribute per user improves data utility, which is an extension of common results in the LDP literature [22], [27], [34]–[36].

On the other hand, in the literature, both SUE and OUE protocols have been extended (and also applied [37], [38]) to longitudinal studies based on the concept of *memoization* [11], [12], i.e., L-SUE and L-OUE, respectively. However, we numerically and experimentally show that combining both protocols provides higher data utility, i.e., starting with OUE and then with SUE (L-OSUE) optimizes data utility rather than using SUE or OUE twice. In addition, we also extended GRR for longitudinal studies (i.e., L-GRR), which provides higher data utility than the other protocols based on unary encoding for attributes with small domain size.

Lastly, in a multidimensional setting with different domain sizes for each attribute, a dynamic selection of longitudinal LDP protocols is preferred. Therefore, we also proposed a new solution named Adaptive LDP for LOngitudinal and Multidimensional FREquency Estimates (ALLOMFREE), which combines all the aforementioned results. More specifically, ALLOMFREE randomly samples a single attribute to send with the whole privacy budget and adaptively selects the optimal protocol, i.e., either L-GRR or L-OSUE. To validate our proposal, we conducted a comprehensive and extensive set of experiments on four real-world open datasets. Under the same privacy guarantee, results show that ALLOMFREE consistently and considerably outperforms the state-of-the-art L-SUE and L-OUE protocols in the quality of the frequency estimates.

**Paper's Outline.** The remainder of this paper is organized as follows. In Section II, we review the privacy notion that we are considering, i.e., LDP and the protocols we further analyze in this paper. In Section III, we extend the analysis of GRR, OUE, and SUE to multidimensional data collections. In Section IV we present the *memoization*-based framework for longitudinal data collections, the extension and analysis

of longitudinal GRR and longitudinal UE-based protocols; the numerical evaluation of their performance, and we present our ALLOMFREE solution. In Section V, we present experimental results, discuss our results and limitations, and review related work. Lastly, in Section VI, we present the concluding remarks and future directions.

## II. THEORETICAL BACKGROUND

In this section, we briefly present the concept of privacy considered in this work, that is, LDP (Subsection II-A), and the LDP protocols we will apply in this paper and their analysis (Subsection II-B).

### A. Local differential privacy (LDP)

Local differential privacy, initially formalized in [10], protects an individual's privacy during the data collection process. A formal definition of LDP is given in the following:

**Definition 1** ($\epsilon$-Local Differential Privacy). *A randomized algorithm $\mathcal{A}$ satisfies $\epsilon$-LDP if, for any pair of input values $v_1, v_2 \in Domain(\mathcal{A})$ and any possible output $y$ of $\mathcal{A}$:*

$$\Pr[\mathcal{A}(v_1) = y] \le e^\epsilon \cdot \Pr[\mathcal{A}(v_2) = y].$$

Similar to the centralized model of DP, LDP also enjoys several important properties, e.g., immunity to post-processing ($F(\mathcal{A})$ is $\epsilon$-LDP for any function $F$) and composability [3]. That is, combining the results from $m$ locally differentially private protocols also satisfies LDP. If these protocols are applied separately in disjointed subsets of the dataset, $\epsilon = max(\epsilon_1$-, ..., $\epsilon_m)$-LDP (parallel composition). On the other hand, if these protocols are sequentially applied to the same dataset, $\epsilon = \sum_{i=1}^m \epsilon_i$-LDP (sequential composition).

### B. LDP protocols

Randomized response (RR), a surveying technique proposed by Warner [39], has been the building block for many LDP protocols. Let $A_j = \{v_1, v_2, ..., v_{k_j}\}$ be a set of $k_j = |A_j|$ values of a given attribute and let $\epsilon$ be the privacy budget, we review three state-of-the-art LDP mechanisms for single-frequency estimation (a.k.a. frequency oracles) that will be used in this paper.

*1) Generalized randomized response (GRR):* The $k$-Ary RR [16] mechanism extends RR to the case of $k_j \ge 2$ and is also referred to as direct encoding [14] or generalized RR (GRR) [27], [40], [41]. Throughout this paper, we use the term GRR for this LDP protocol. Given a value $v \in \mathcal{D}_j$, *GRR(v)* outputs the true value with probability $p$, and any other value $v' \in \mathcal{D}_j$ such that $v' \neq v$ with probability $1-p$. More formally, the perturbation function is defined as:

$$\forall y \in \mathcal{D}_j \; \Pr[\mathcal{A}_{GRR(\epsilon)}(v) = y] = \begin{cases} p = \frac{e^\epsilon}{e^\epsilon + k_j - 1}, & \text{if } y = v \\ q = \frac{1}{e^\epsilon + k_j - 1}, & \text{if } y \neq v. \end{cases}$$

This satisfies $\epsilon$-LDP since $\frac{p}{q} = e^\epsilon$. To estimate the frequency $f(v_i)$ that a value $v_i$ occurs for $i \in [1, k_j]$, one calculates [14]:

$$\hat{f}(v_i) = \frac{N_i - nq}{n(p - q)}, \qquad (1)$$

in which $N_i$ is the number of times the value $v_i$ has been reported and $n$ is the total number of users. In [14], it is shown that $\hat{f}(v_i)$ is an unbiased estimation of the true frequency $f(v_i)$, and the variance of this estimation is $Var[\hat{f}(v_i)] = \frac{q(1-q)}{n(p-q)^2} + \frac{f(v_i)(1-p-q)}{n(p-q)}$. In the case of small $f(v_i) \sim 0$, this variance is dominated by the first term, which gives the *approximate* variance as [14]:

$$Var^*[\hat{f}(v_i)] = \frac{q(1-q)}{n(p-q)^2}. \qquad (2)$$

Replacing $p = \frac{e^\epsilon}{e^\epsilon + k_j - 1}$ and $q = \frac{1}{e^\epsilon + k_j - 1}$ into Eq. (2), the GRR variance is calculated as:

$$Var^*[\hat{f}_{GRR}(v_i)] = \frac{e^\epsilon + k_j - 2}{n(e^\epsilon - 1)^2}. \qquad (3)$$

*2) Unary encoding-based:* Protocols based on unary encoding (UE) consist of transforming a value $v$ into a binary representation of it. So, first, for a given value $v$, $B = Encode(v)$, where $B = [0, 0, ..., 1, 0, ...0]$, a $k_j$-bit array where only the $v$-th position is set to one. Next, the bits from $B$ are flipped, depending on parameters $p$ and $q$, to generate a sanitized vector $B'$, in which:

$$Pr[B'[i] = 1] = \begin{cases} p, & \text{if } B[i] = 1 \\ q, & \text{if } B[i] = 0. \end{cases}$$

The proof that UE-based protocols satisfy $\epsilon$-LDP for

$$\epsilon = ln\left(\frac{p(1-q)}{(1-p)q}\right), \qquad (4)$$

is known in the literature and can be found in [11], [14]. In [14] the authors presents two ways for selecting probabilities $p$ and $q$, which determines the protocol variance. One well-known UE-based protocol is the Basic One-time RAPPOR [11], referred to as symmetric UE (SUE), which selects $p = \frac{e^{\epsilon/2}}{e^{\epsilon/2}+1}$ and $q = \frac{1}{e^{\epsilon/2}+1}$, where $p + q = 1$ (symmetric). The estimated frequency $\hat{f}(v_i)$ that a value $v_i$ occurs for $i \in [1, k_j]$ is also calculated using Eq. (1). Replacing $p = \frac{e^{\epsilon/2}}{e^{\epsilon/2}+1}$ and $q = \frac{1}{e^{\epsilon/2}+1}$ into Eq. (2), the SUE variance is calculated as [11]:

$$Var^*[\hat{f}_{SUE}(v_i)] = \frac{e^{\epsilon/2}}{n(e^{\epsilon/2} - 1)^2}. \qquad (5)$$

Moreover, rather than selecting $p$ and $q$ to be symmetric, Wang et al. [14] proposed optimized UE (OUE), which selects parameters $p = \frac{1}{2}$ and $q = \frac{1}{e^\epsilon+1}$ that minimize the variance of UE-based protocols while still satisfying $\epsilon$-LDP. Similarly, the estimation method used in Eq. (1) equally applies to OUE. Replacing $p = \frac{1}{2}$ and $q = \frac{1}{e^\epsilon+1}$ into Eq. (2), the OUE variance is calculated as [14]:

$$Var^*[\hat{f}_{OUE}(v_i)] = \frac{4e^\epsilon}{n(e^\epsilon - 1)^2}. \qquad (6)$$

## III. MULTIDIMENSIONAL FREQUENCY ESTIMATES WITH LDP

In the literature, there are few works for collecting multidimensional data with LDP based on random sampling (i.e., dividing users in groups) [14], [30]–[33], [36]. This technique reduces both dimensionality and communication costs, which will also be the focus of this paper. Let $d \geq 2$ be the total number of attributes, $\mathbf{k} = [k_1, k_2, ..., k_d]$ be the domain size of each attribute, $n$ be the number of users, and $\epsilon$ be the privacy budget. An intuitive solution (*Spl*) is splitting the privacy budget, i.e., assigning $\epsilon/d$ for each attribute. The other solution (*Smp*) is based on uniformly sampling (without replacement) only $r$ attribute(s) out of $d$ possible ones, i.e., assigning $\epsilon/r$ per attribute. Notice that both solutions satisfy $\epsilon$-LDP according to the sequential composition theorem [3].

For the first case, *Spl*, the variances ($\sigma_1^2$) of GRR, SUE, and OUE are, respectively:

$$\sigma_{1,GRR}^2 = \frac{e^{\epsilon/d} + k_j - 2}{n(e^{\epsilon/d} - 1)^2},$$
$$\sigma_{1,SUE}^2 = \frac{e^{\epsilon/2d}}{n(e^{\epsilon/2d} - 1)^2}, \qquad (7)$$
$$\sigma_{1,OUE}^2 = \frac{4e^{\epsilon/d}}{n(e^{\epsilon/d} - 1)^2}.$$

For the second case, *Smp*, the number of users per attribute is reduced to $nr/d$. Thus, the variances ($\sigma_2^2$) of GRR, SUE, and OUE are, respectively:

$$\sigma_{2,GRR}^2 = \frac{d(e^{\epsilon/r} + k_j - 2)}{nr(e^{\epsilon/r} - 1)^2},$$
$$\sigma_{2,SUE}^2 = \frac{d(e^{\epsilon/2r})}{nr(e^{\epsilon/2r} - 1)^2}, \qquad (8)$$
$$\sigma_{2,OUE}^2 = \frac{d(4e^{\epsilon/r})}{nr(e^{\epsilon/r} - 1)^2}.$$

Notice that if $r = d$ in Eq. (8), one achieves Eq. (7). Practically, the objective is reduced to finding $r$, which minimizes $\sigma_2^2$ for each protocol. This way, to find the optimal $r$ for each protocol, we first multiply each $\sigma_2^2$ in Eq. (8) by $\epsilon$. Without loss of generality, minimizing $\sigma_{2,GRR}^2$, $\sigma_{2,SUE}^2$, and $\sigma_{2,OUE}^2$ is equivalent to minimizing $\frac{\epsilon e^{\epsilon/r}}{r(e^{\epsilon/r}-1)^2}$, $\frac{\epsilon e^{\epsilon/2r}}{r(e^{\epsilon/2r}-1)^2}$, and $\frac{\epsilon e^{\epsilon/r}}{r(e^{\epsilon/r}-1)^2}$, respectively. Hence, let $x = r/\epsilon$ be the independent variable, $\sigma_{2,GRR}^2$ and $\sigma_{2,OUE}^2$ can be rewritten as $y_1 = \frac{1}{x} \cdot \frac{e^{1/x}}{(e^{1/x}-1)^2}$, and $\sigma_{2,SUE}^2$ can be rewritten as $y_2 = \frac{1}{x} \cdot \frac{e^{1/2x}}{(e^{1/2x}-1)^2}$ as functions over $x$. It is not hard to prove that both $y_1$ and $y_2$ are increasing functions w.r.t. $x$ and, hence, we have a minimum and optimal when $r = 1$ (a single attribute per user) for all three protocols. We highlight that this is a common result in the LDP literature obtained for different protocols and contexts [14], [22], [30], [31], [33]–[35], [42].

**Therefore, in this paper, we adopt the multidimensional setting *Smp* with $r = 1$. In this setting, users tell the data**

collector which attribute was sampled, and its perturbed value ensuring $\epsilon$-LDP by applying either GRR or UE-based protocols; the data analyst server would not receive any information about the remaining $d - 1$ attributes.

## IV. LONGITUDINAL FREQUENCY ESTIMATES WITH LDP

In this section, we present the *memoization*-based framework for longitudinal data collections (Subsection IV-A). Next, we present the analysis of longitudinal GRR (Subsection IV-B) and longitudinal UE-based protocols (Subsection IV-C). Lastly, we evaluate numerically the extended longitudinal protocols (Subsection IV-D) and we propose our ALLOMFREE solution (Subsection IV-E).

### A. Memoization-based data collection with LDP

In the literature, many works study how to collect and analyze categorical data longitudinally based on *memoization* [11], [12], [34]. The key idea behind memoization is using two sanitization processes. The first round ($RR_1$) replaces the real value $B$ with a sanitized one $B'$ with a higher epsilon ($\epsilon_\infty$). Whenever one intends to report $B$, $B'$ shall be reused to produce other sanitized versions $B''$ with lower epsilon values. Notice that the second sanitization ($RR_2$) is a *must* to avoid 'averaging attacks', in which adversaries can reconstruct the true value from multiple sanitized versions of it. This technique allows achieving privacy over time with an upper bound value of $\epsilon_\infty$-LDP.

Let $A_j = \{v_1, v_2, ..., v_{k_j}\}$ be a set of $k_j = |A_j|$ values of a given attribute and let $\epsilon$ be the privacy budget. In this paper, for both $RR_1$ and $RR_2$ steps, we will apply either GRR, SUE, or OUE. The unbiased estimator in Eq. (1) for the frequency $f(v_i)$ of each value $v_i$ for $i \in [1, k_j]$ is now extended to:

$$\hat{f}_L(v_i) = \frac{N_i - nq_1(p_2 - q_2) - nq_2}{n(p_1 - q_1)(p_2 - q_2)}, \qquad (9)$$

in which $N_i$ is the number of times the value $v_i$ has been reported, $n$ is the total number of users, $p_1$ and $q_1$ are the parameters used by an LDP protocol for $RR_1$, and $p_2$ and $q_2$ are the parameters used by an LDP protocol for $RR_2$.

**Theorem 1.** *The estimation result $\hat{f}_L(v_i)$ in Eq.* (9) *is an unbiased estimation of $f(v_i)$ for any value $v_i \in \mathcal{D}_j$.*

*Proof 1*

$$E[\hat{f}_L(v_i)] = E\left[\frac{N_i - nq_1(p_2 - q_2) + nq_2}{n(p_1 - q_1)(p_2 - q_2)}\right]$$
$$= \frac{E[Ni]}{n(p_1 - q_1)(p_2 - q_2)} - \frac{q_1(p_2 - q_2) - q_2}{(p_1 - q_1)(p_2 - q_2)}.$$

Let us focus on

$$E[N_i] = nf(v_i)(p_1 p_2 + q_2(1 - p_1))$$
$$+ n(1 - f(v_i))(p_2 q_1 + q_2(1 - q_1)).$$

Thus,

$$E[\hat{f}_L(v_i)] = f(v_i).$$

**Theorem 2.** *The variance of the estimation in Eq.* (9) *is:*

$$Var[\hat{f}_L(v_i)] = \frac{\gamma(1 - \gamma)}{n(p_1 - q_1)^2(p_2 - q_2)^2}, \text{ where}$$
$$\gamma = f(v_i)(2p_1 p_2 - 2p_1 q_2 + 2q_2 - 1) + p_2 q_1 + q_2(1 - q_1). \tag{10}$$

*Proof 2*

Thanks to Eq. (9) we have

$$Var[\hat{f}_L(v_i)] = \frac{Var(N_i)}{n^2(p_1 - q_1)^2(p_2 - q_2)^2}.$$

Since $N_i$ is the number of times the value $v_i$ is observed, it can be defined as $N_i = \sum_{z=1}^{n} X_z$ where $X_z$ is equal to 1 if the user $z$, $1 \leq z \leq n$ reports value $v_i$, and 0 otherwise. We thus have $Var(N_i) = \sum_{z=1}^{n} Var(X_z) = nVar(X)$. Since all the users are independent,

$$P(X = 1) = P(X^2 = 1) = f(v_i)(2p_1 p_2 - 2p_1 q_2 + 2q_2 - 1)$$
$$+ p_2 q_1 + q_2(1 - q_1) = \gamma.$$

We thus have $Var(X) = \gamma - \gamma^2 = \gamma(1 - \gamma)$ and, finally,

$$Var[\hat{f}_L(v_i)] = \frac{\gamma(1 - \gamma)}{n(p_1 - q_1)^2(p_2 - q_2)^2}.$$

In this work, we will use the *approximate variance*, in which $f(v_i) = 0$ in Eq. (10), which gives:

$$Var^*[\hat{f}_L(v_i)] =$$
$$\frac{(p_2 q_1 - q_2(q_1 - 1))(-p_2 q_1 + q_2(q_1 - 1) + 1)}{n(p_1 - q_1)^2(p_2 - q_2)^2}. \tag{11}$$

### B. Longitudinal GRR (L-GRR): definition and $\epsilon$-LDP study

Let $V = \{v_1, v_2, ..., v_{k_j}\}$ be a set of $k_j$ values of a given attribute and let $v_i$ be the real value. We now describe an extension of GRR for longitudinal studies; we refer to this protocol as L-GRR for the rest of this paper. First, $Encode(v_i) = v_i$ (direct encoding). Next, there are two rounds of sanitization, $RR_1$ and $RR_2$ applying GRR, described in the following.

1) $RR_1[GRR]$: Memoize a value $B'$ such that

$$B' = \begin{cases} v_i, & \text{with probability } p_1, \\ v_{k \neq v_i}, & \text{with probability } q_1 = \frac{1 - p_1}{k_j - 1}, \end{cases}$$

in which $p_1$ and $q_1$ control the level of longitudinal $\epsilon_\infty$-LDP. The value $B'$ shall be reused as the basis for all future reports on the real value $v_i$.

2) $RR_2[GRR]$: Generate a reporting $B''$ such that

$$B'' = \begin{cases} B', & \text{with probability } p_2, \\ v_{k \neq B'}, & \text{with probability } q_2 = \frac{1 - p_2}{k_j - 1}, \end{cases}$$

in which $B''$ is the report to be sent to the server.

Visually, Fig. 1 illustrates the probability tree of the L-GRR protocol. In the first round of sanitization, $RR_1$, our proposed L-GRR applies GRR with $p_1 = Pr[B' = v_i | B = v_i] = \frac{e^{\epsilon_\infty}}{e^{\epsilon_\infty} + k_j - 1}$ and $q_1 = Pr[B' = v_i | B = v_{k \neq i}] =$
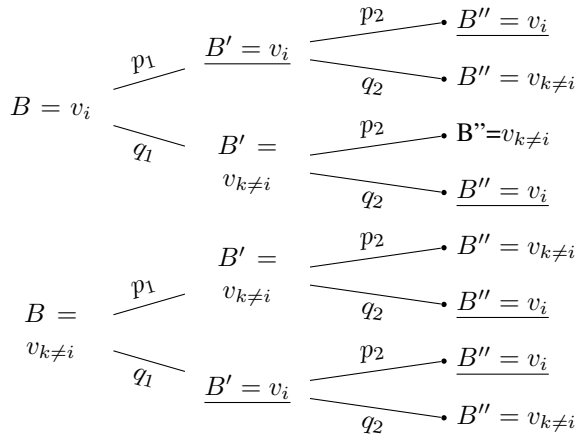
Fig. 1: Probability tree for two rounds of sanitization using GRR (L-GRR).



Fig. 2: Probability tree for two rounds of sanitization using UE (L-UE).

$\frac{1-p_1}{k_j-1} = \frac{1}{e^{\epsilon_\infty}+k_j-1}$ (underlined in the middle of Fig. 1), where $k_j = |A_j|$. As discussed in Subsection II-B1, this *permanent memoization* satisfies $\epsilon_\infty$-LDP since $\frac{p_1}{q_1} = e^{\epsilon_\infty}$, which is the upper bound.

On the other hand, with a single collection of data, the attacker's knowledge of $v_i$ comes only from $B''$, which is generated using two randomization steps with GRR. This provides a higher level of privacy protection [11]. From Fig. 1, we can obtain the following conditional probabilities:

$$\Pr[B''|B] = \begin{cases} \Pr[B'' = v_i|B = v_i] = p_1p_2 + q_1q_2 \\ \Pr[B'' = v_{k\neq i}|B = v_i] = p_1q_2 + q_1p_2 \\ \Pr[B'' = v_i|B = v_{k\neq i}] = p_1q_2 + q_1p_2 \\ \Pr[B'' = v_{k\neq i}|B = v_{k\neq i}] = p_1p_2 + q_1q_2 \end{cases}$$

Let $p_s = \Pr[B'' = v_i|B = v_i]$ and $q_s = \Pr[B'' = v_i|B = v_{k\neq i}]$ (underlined in far right of Fig. 1), with the second round of sanitization, $RR_2[GRR]$, our proposed L-GRR protocol satisfies $\epsilon_1$-LDP since $\frac{p_s}{q_s} = e^{\epsilon_1}$. Notice that $\epsilon_1$ corresponds to a single report (lower bound) and its extension to infinity reports is limited by $\epsilon_\infty$ (upper bound) since $RR_2[GRR]$ uses as input the output of $RR_1[GRR]$. More specifically, the calculus of $\epsilon_1$ for L-GRR is:

$$\epsilon_1 = \ln\left(\frac{p_1p_2 + q_1q_2}{p_1q_2 + q_1p_2}\right), \tag{12}$$

in which $p_1 = \frac{e^{\epsilon_\infty}}{e^{\epsilon_\infty}+k_j-1}$, $q_1 = \frac{1-p_1}{k_j-1}$, and both $p_2$ and $q_2$ are selectable according with $\epsilon_\infty$, $\epsilon_1$, and $k_j$, calculated as:

$$p_2 = \frac{e^{\epsilon_1+\epsilon_\infty}-1}{-k_je^{\epsilon_1}+(k_j-1)e^{\epsilon_\infty}+e^{\epsilon_1}+e^{\epsilon_1+\epsilon_\infty}-1},$$
$$q_2 = \frac{1-p_2}{k_j-1}. \tag{13}$$

The estimated frequency $\hat{f}_L(v_i)$ that a value $v_i$ occurs for $i \in [1, k_j]$ is calculated using Eq. (9). Lastly, one can calculate the L-GRR approximate variance by replacing the resulting $p_1, q_1, p_2, q_2$ parameters into Eq. (11).
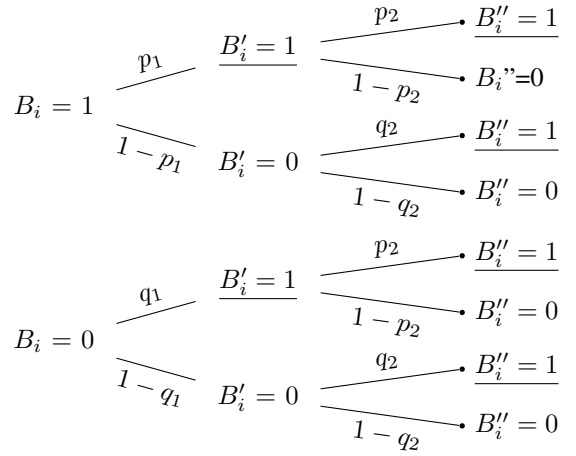
## C. Longitudinal UE (L-UE): definition and $\epsilon$-LDP study

We now describe UE-based protocols for longitudinal studies; we refer to this protocol as L-UE for the rest of this paper. Let $V = \{v_1, v_2, ..., v_{k_j}\}$ be a set of $k_j$ values of a given attribute and let $v_i$ be the real value. First, $Encode(v_i) = B$ (unary encoding), where $B = [0, 0, ..., 1, 0, ...0]$, a $k_j$-bit array where only the $v$-th position is set to one. Next, there are two rounds of sanitization, $RR_1$ and $RR_2$ applying UE-based protocols, described in the following.

1) $RR_1[UE]$: For each bit $i$, $1 \leq i \leq k_j$ in $B$, memoize a value $B'$ such that
$$P(B_i' = 1) = \begin{cases} p_1, & \text{if } B_i = 1 \text{ and} \\ q_1, & \text{if } B_i = 0, \end{cases}$$
in which $p_1$ and $q_1$ control the level of longitudinal $\epsilon_\infty$-LDP. The value $B'$ shall be reused as the basis for all future reports on the real value $v_i$.

2) $RR_2[UE]$: For each bit $i$, $1 \leq i \leq k_j$ in $B'$, generate a reporting $B''$ that
$$P(B_i'' = 1) = \begin{cases} p_2, & \text{if } B_i' = 1 \text{ and} \\ q_2, & \text{if } B_i' = 0, \end{cases}$$
in which $B''$ is the report to be sent to the server.

Visually, Fig. 2 illustrates the probability tree of the L-UE protocol. **One natural question emerges: how to select the parameters $\{p_1, q_1, p_2, q_2\}$ in order to optimize the utility of this L-UE protocol?** One can see $RR_1[UE]$ as a *permanent* sanitization and $RR_2[UE]$ as a 'small' perturbation to avoid averaging attacks and keep privacy over time.

Based on SUE and OUE, we are then left with four options: two popular solutions that strictly use only OUE or SUE parameters in both sanitization steps and two proposed settings that combine both OUE and SUE. These four L-UE protocols are summarized below:

  I  both sanitizations with OUE (L-OUE);
  II  both sanitizations with SUE (L-SUE);
  III  starting with OUE and then with SUE (L-OSUE);
  IV  starting with SUE and then with OUE (L-SOUE);

in which, L-SUE is the well-known Basic-RAPPOR protocol [11], L-OUE is the state-of-the-art OUE protocol [14] with memoization, and both L-OSUE and L-SOUE are proposed in this paper.

As presented in [14], the OUE variance in Eq. (6) is smaller than the SUE variance in Eq. (5) and, therefore, the former can provide higher utility than the latter for $RR_1$. On the other hand, we argue that OUE might be too strict for $RR_2$ since the parameter $p_2 = 1/2$ is constant. Thus, we hypothesize that option III (i.e., L-OSUE) is the most suitable one. Without loss of generality, **the following analyses are done only for L-OSUE**, which can be easily extended to any of the other combinations.

In the first round of sanitization, $RR_1$, our solution L-OSUE applies OUE with $p_1 = Pr[B_i' = 1 | B_i = 1] = \frac{1}{2}$ and $q_1 = Pr[B_i' = 1 | B_i = 0] = \frac{1}{e^{\epsilon_\infty} + 1}$ (underlined in the middle of Fig. 2). As discussed in Section II-B2, this *permanent* memoization satisfies $\epsilon_\infty$-LDP since $\frac{p_1(1-q_1)}{(1-p_1)q_1} = e^{\epsilon_\infty}$, which is the upper bound.

Following the same development as for L-GRR, on the other hand, with a single collection of data, the attacker's knowledge of $B = Encode(v)$ comes only from $B''$, which is generated using two randomization steps with OUE and SUE, respectively. This provides a higher level of privacy protection [11]. From Fig. 2, we can obtain the following conditional probabilities according to each bit $i \in [1, k_j]$:

$$Pr[B_i''|B_i] =$$

$$\begin{cases} Pr[B_i'' = 1|B_i = 1] = p_1 p_2 + (1 - p_1)q_2 \\ Pr[B_i'' = 0|B_i = 1] = p_1(1 - p_2) + (1 - p_1)(1 - q_2) \\ Pr[B_i'' = 1|B_i = 0] = q_1 p_2 + (1 - q_1)q_2 \\ Pr[B_i'' = 0|B_i = 0] = q_1(1 - p_2) + (1 - q_1)(1 - q_2) \end{cases}$$

Let $p_s = Pr[B_i'' = 1|B_i = 1]$ and $q_s = Pr[B_i'' = 1|B_i = 0]$ (underlined in far right of Fig. 2), with the second round of sanitization, $RR_2[SUE]$, our proposed L-OSUE protocol satisfies $\epsilon_1$-LDP since $\frac{p_s(1-q_s)}{(1-p_s)q_s} = e^{\epsilon_1}$. Notice that $\epsilon_1$ corresponds to a single report (lower bound) and its extension to infinity reports is limited by $\epsilon_\infty$ (upper bound) since $RR_2[SUE]$ uses as input the output of $RR_1[OUE]$. More specifically, the calculus of $\epsilon_1$ for L-OSUE (or L-UE protocols in general) is:

$$\epsilon_1 = \ln\left(\frac{(p_1 p_2 - q_2(p_1 - 1))(p_2 q_1 - q_2(q_1 - 1) - 1)}{(p_2 q_1 - q_2(q_1 - 1))(p_1 p_2 - q_2(p_1 - 1) - 1)}\right),$$

(14)

in which, for L-OSUE, we have $p_1 = \frac{1}{2}$, $q_1 = \frac{1}{e^{\epsilon_\infty} + 1}$, and both $p_2$ and $q_2$ are symmetric ($p_2 + q_2 = 1$) and selectable according to $\epsilon_\infty$ and $\epsilon_1$, calculated as:

$$p_2 = \frac{1 - e^{\epsilon_1 + \epsilon_\infty}}{e^{\epsilon_1} - e^{\epsilon_\infty} - e^{\epsilon_1 + \epsilon_\infty} + 1},$$

(15)

$$q_2 = 1 - p_2.$$

Similarly, the estimated frequency $\hat{f}_L(v_i)$ that a value $v_i$ occurs for $i \in [1, k_j]$ is calculated using Eq. (9). Lastly, one can calculate the L-OSUE (or L-UE protocols in general) approximate variance by replacing the resulting $p_1, q_1, p_2, q_2$ parameters into Eq. (11).

### D. Numerical evaluation of L-GRR and L-UE protocols

In this subsection, we evaluate numerically the approximate variance of all developed longitudinal protocols, namely, L-GRR and the four UE-based options namely L-OUE, L-SUE, L-OSUE, and L-SOUE, respectively. As aforementioned, once defined both $\epsilon_\infty$ and $\epsilon_1$ privacy guarantees, one can obtain the parameters $p_1$ and $q_1$ depending on $\epsilon_\infty$, and the parameters $p_2$ and $q_2$ depending on both $\epsilon_\infty$ and $\epsilon_1$ (and the domain size $k_j$ for L-GRR) as given in Eq. (13) for L-GRR and in Eq. (15) for L-OSUE.

Next, once computed the parameters $\{p_1, q_1, p_2, q_2\}$, one can calculate the approximate variance with Eq. (11) for each protocol. In other words, following our proposal, one has to set both the upper ($\epsilon_\infty$) and lower ($\epsilon_1$) bounds of the privacy guarantees. For example, let $\epsilon_\infty = 2$, one might want that the first $\epsilon_1$-LDP report to have high privacy such as $\epsilon_1 = 0.1$, i.e., $\epsilon_1 = 0.05\epsilon_\infty$ (**we will use this probability notation to set up the privacy guarantees**).

Table I exhibits numerical values of the approximate variance using Eq. (11) for all longitudinal protocols with $n = 10000$, $\epsilon_\infty = [0.5, 1.0, 2.0, 4.0]$ (as in [14]), and $\epsilon_1 = \{0.6\epsilon_\infty, 0.5\epsilon_\infty, 0.4\epsilon_\infty, 0.3\epsilon_\infty, 0.2\epsilon_\infty, 0.1\epsilon_\infty\}$. For values of $\epsilon_1$ higher than $0.6\epsilon_\infty$, neither L-OUE nor L-SOUE could satisfy some values of $\epsilon_1$ because of the constant $p_2 = 1/2$ in $RR_2$. Yet, it is not desirable to have higher values of $\epsilon_1$ and, thus, we did not consider values above $0.6\epsilon_\infty$ in our analysis. Besides, Table II exhibits numerical values for non-longitudinal GRR, OUE, and SUE protocols, which allows evaluating how utility degrades with a second step of sanitization.

**From Table I, one can notice that L-GRR presents the smallest variance values for binary attributes (i.e., when $k_j = 2$).** On the other hand, L-GRR is also the most sensitive to change in privacy parameters $\epsilon_\infty$ and $\epsilon_1$ when $k_j$ is large, which leads to much higher variance than when using a non-longitudinal GRR in Table II. Similar to non-longitudinal GRR, this increase in the variance is due to the number of values $k_j$, which decreases the probability $p$ of reporting the true value. With two rounds of sanitization, it further deteriorates the accuracy of the L-GRR protocol getting to extremely high values, e.g., see L-GRR($k_j = 2^{10}$). Interestingly, when $k_j = 2$ in Table I, the variance of L-GRR with $\epsilon_1 = 0.5\epsilon_\infty$ is a lagged version of the variance values given by the non-longitudinal GRR in Table II. This effect is also observed for both L-SUE (cf. SUE in Table II) and L-OSUE (cf. OUE in Table II) protocols, which use symmetric probabilities on $RR_2$ (i.e., $p_2 + q_2 = 1$). We highlighted these values in **bold font**. However, for L-GRR, this is not true for other values of $k_j$, whose further analysis is beyond the scope of this paper.

On the other hand, L-UE protocols avoid having a variance that depends on $k_j$ by encoding the value into the unary representation, which results in a constant variance no matter

| Privacy Guarantees | | L-GRR | | | L-UE | | | |
|---|---|---|---|---|---|---|---|---|
| | | $k_j = 2$ | $k_j = 32$ | $k_j = 2^{10}$ | L-OSUE | L-SUE | L-SOUE | L-OUE |
| $\epsilon_1 = 0.6\epsilon_\infty$ | $\epsilon_\infty = 0.5, \epsilon_1 = 0.30$ | 0.001103 | 0.980969 | 26706 | 0.004411 | 0.004436 | 0.005306 | 0.005549 |
| | $\epsilon_\infty = 1.0, \epsilon_1 = 0.60$ | 0.000270 | 0.125036 | 3153 | 0.001078 | 0.001103 | 0.001234 | 0.001347 |
| | $\epsilon_\infty = 2.0, \epsilon_1 = 1.20$ | 0.000062 | 0.006327 | 117 | 0.000247 | 0.000270 | 0.000264 | 0.000310 |
| | $\epsilon_\infty = 4.0, \epsilon_1 = 2.40$ | 0.000011 | 0.000078 | 0.25903 | 0.000044 | 0.000062 | 0.000045 | 0.000057 |
| $\epsilon_1 = 0.5\epsilon_\infty$ | $\epsilon_\infty = 0.5, \epsilon_1 = 0.25$ | 0.001592 | 2.088372 | 60218 | 0.006367 | 0.006392 | 0.007336 | 0.007611 |
| | $\epsilon_\infty = 1.0, \epsilon_1 = 0.50$ | **0.000392** | 0.268074 | 7198 | **0.001567** | **0.001592** | 0.001740 | 0.001872 |
| | $\epsilon_\infty = 2.0, \epsilon_1 = 1.00$ | **0.000092** | 0.013926 | 281 | **0.000368** | **0.000392** | 0.000389 | 0.000447 |
| | $\epsilon_\infty = 4.0, \epsilon_1 = 2.00$ | **0.000018** | 0.000188 | 0.74088 | **0.000072** | **0.000092** | 0.000073 | 0.000092 |
| $\epsilon_1 = 0.4\epsilon_\infty$ | $\epsilon_\infty = 0.5, \epsilon_1 = 0.20$ | 0.002492 | 4.530779 | 135874 | 0.009967 | 0.009992 | 0.011012 | 0.011324 |
| | $\epsilon_\infty = 1.0, \epsilon_1 = 0.40$ | 0.000617 | 0.586823 | 16443 | 0.002467 | 0.002492 | 0.002658 | 0.002812 |
| | $\epsilon_\infty = 2.0, \epsilon_1 = 0.80$ | 0.000148 | 0.031552 | 673 | 0.000593 | 0.000617 | 0.000617 | 0.000690 |
| | $\epsilon_\infty = 4.0, \epsilon_1 = 1.60$ | 0.000032 | 0.000484 | 2.12772 | 0.000127 | 0.000148 | 0.000128 | 0.000156 |
| $\epsilon_1 = 0.3\epsilon_\infty$ | $\epsilon_\infty = 0.5, \epsilon_1 = 0.15$ | 0.004436 | 10 | 329836 | 0.017744 | 0.017769 | 0.018863 | 0.019214 |
| | $\epsilon_\infty = 1.0, \epsilon_1 = 0.30$ | 0.001103 | 1.398568 | 40412 | 0.004411 | 0.004436 | 0.004620 | 0.004799 |
| | $\epsilon_\infty = 1.0, \epsilon_1 = 0.60$ | 0.000270 | 0.078202 | 1737 | 0.001078 | 0.001103 | 0.001106 | 0.001198 |
| | $\epsilon_\infty = 2.0, \epsilon_1 = 1.20$ | 0.000062 | 0.001389 | 6 | 0.000247 | 0.000270 | 0.000248 | 0.000291 |
| $\epsilon_1 = 0.2\epsilon_\infty$ | $\epsilon_\infty = 0.5, \epsilon_1 = 0.10$ | 0.009992 | 30 | 972656 | 0.039967 | 0.039992 | 0.041148 | 0.041536 |
| | $\epsilon_\infty = 1.0, \epsilon_1 = 0.20$ | 0.002492 | 4.080052 | 120651 | 0.009967 | 0.009992 | 0.010190 | 0.010394 |
| | $\epsilon_\infty = 2.0, \epsilon_1 = 0.40$ | 0.000617 | 0.237925 | 5443 | 0.002467 | 0.002492 | 0.002498 | 0.002610 |
| | $\epsilon_\infty = 4.0, \epsilon_1 = 0.80$ | 0.000148 | 0.004939 | 24 | 0.000593 | 0.000617 | 0.000595 | 0.000659 |
| $\epsilon_1 = 0.1\epsilon_\infty$ | $\epsilon_\infty = 0.5, \epsilon_1 = 0.05$ | 0.039992 | 154 | 4941829 | 0.159967 | 0.159992 | 0.161191 | 0.161608 |
| | $\epsilon_\infty = 1.0, \epsilon_1 = 0.10$ | 0.009992 | 20 | 620584 | 0.039967 | 0.039992 | 0.040201 | 0.040424 |
| | $\epsilon_\infty = 2.0, \epsilon_1 = 0.20$ | 0.002492 | 1.255550 | 29356 | 0.009967 | 0.009992 | 0.010000 | 0.010130 |
| | $\epsilon_\infty = 4.0, \epsilon_1 = 0.40$ | 0.000617 | 0.030494 | 156 | 0.002467 | 0.002492 | 0.002469 | 0.002560 |

TABLE I: Numerical values of Eq. (11) (i.e., $Var^*[\hat{f}_L(v_i)]$) for L-GRR and L-UE protocols with different $\epsilon_\infty$ and $\epsilon_1$ privacy guarantees, following $\epsilon_1 = \{0.6\epsilon_\infty, 0.5\epsilon_\infty, 0.4\epsilon_\infty, 0.3\epsilon_\infty, 0.2\epsilon_\infty, 0.1\epsilon_\infty\}$, respectively.

| $\epsilon_\infty$ | GRR($k_j = 2$) | GRR($k_j = 32$) | GRR($k_j = 2^{10}$) | OUE | SUE |
|---|---|---|---|---|---|
| $\epsilon_\infty = 0.5$ | **0.000392** | 0.007520 | 0.243240 | **0.001567** | **0.001592** |
| $\epsilon_\infty = 1.0$ | **0.000092** | 0.001108 | 0.034707 | **0.000368** | **0.000392** |
| $\epsilon_\infty = 2.0$ | **0.000018** | 0.000092 | 0.002522 | **0.000072** | **0.000092** |
| $\epsilon_\infty = 4.0$ | 0.000002 | 0.000003 | 0.000037 | 0.000008 | 0.000018 |

TABLE II: Numerical values of $Var^*[\hat{f}(v_i)]$ for the non-longitudinal GRR, OUE, and SUE protocols with different $\epsilon_\infty$ privacy guarantees.

the size of the attribute. To complement the results of Table I, Fig. 3 illustrates numerical values of the approximate variance for L-UE protocols with $\epsilon_1 = \{0.3\epsilon_\infty, 0.6\epsilon_\infty\}$. With the four options I-IV analyzed, on high privacy regimes, L-OSUE and L-SUE have similar performance while *always* favoring the proposed L-OSUE one. On lower privacy regimes, our proposed protocols L-SOUE and L-OSUE have similar performance, which outperform both L-OUE and L-SUE protocols. As shown in our experiments, the L-OUE protocol has the worst performance among the four options analyzed, with the exception of high values for $\epsilon_\infty$ (see the plot on the bottom of Fig. 3), when it has performance superior or similar to L-SUE. Indeed, for L-OUE, selecting $p_2 = 1/2$ for the second sanitization step is too strict, which results in higher variance value. **Therefore, by comparing the approximate variances, the best option for L-UE protocols, in terms of utility, is starting with OUE and then with SUE as we propose in this paper, i.e., L-OSUE.**

### E. The ALLOMFREE algorithm

Let $A = \{A_1, A_2, ..., A_d\}$ be a set of $d$ attributes with domain size $\mathbf{k} = [k_1, k_2, ..., k_d]$, $\mathbb{A} = \{L\text{-}GRR, L\text{-}OSUE\}$ be a set of optimal longitudinal LDP protocols, and $\epsilon_\infty$ and $\epsilon_1$ be the longitudinal and *single-report* privacy guarantees, respectively. Each user $u_i$, for $1 \leq i \leq n$, holds a tuple $\mathbf{v}^{(i)} = (v_1^{(i)}, v_2^{(i)}, ..., v_d^{(i)})$, i.e., a private value per attribute. From now on, we will simply omit the index notation $\mathbf{v}^{(i)}$ and use $\mathbf{v}$ in the analysis as we focus on one arbitrary user $u_i$ here. For each attribute $j \in [1, d]$ (we slightly abuse the notation and use $j$ for $A_j$) at time $t \in [1, \tau]$, the aggregator aims to estimate the frequencies of each value $v \in A_j$.

In a multidimensional setting with different domain sizes for each attribute, a dynamic selection of longitudinal LDP protocols is preferred. As mentioned in Section III, we propose that each user randomly sample $r = Uniform(1, 2, ..., d)$ to select a single attribute $A_r$. Given $k_r$ (the domain size), $\epsilon_\infty$, and $\epsilon_1$, one calculates the parameters $fp_{GRR} = \{p_1, q_1, p_2, q_2\}$ and $fp_{UE} = \{p_1, q_1, p_2, q_2\}$, for L-GRR and L-OSUE, respectively (cf. Eq. (13) and Eq. (15)). Next, with $fp_{GRR}$ and $fp_{UE}$, one calculates the approximate variances $Var^*[\hat{f}_{L_{(L\text{-}GRR)}}]$ for L-GRR and $Var^*[\hat{f}_{L_{(L\text{-}OSUE)}}]$ for L-OSUE with Eq. (11). Lastly, to select L-GRR as the local randomizer, we are then left to evaluate if $Var^*[\hat{f}_{L_{(L\text{-}GRR)}}] \leq Var^*[\hat{f}_{L_{(L\text{-}OSUE)}}]$.

We call our solution <u>A</u>daptive <u>L</u>DP for <u>LO</u>ngitudinal
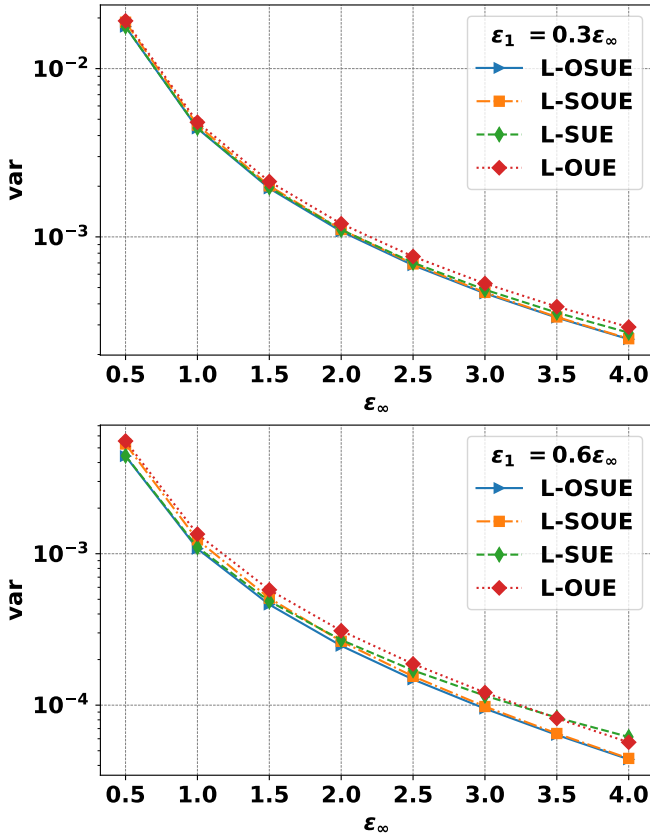
Fig. 3: Numerical values of $Var^*[\hat{f}_L(v_i)]$ for L-UE protocols with $\epsilon_1 = 0.3 \cdot \epsilon_\infty$ (plot on the top) and with $\epsilon_1 = 0.6 \cdot \epsilon_\infty$ (plot on the bottom).

---

**Algorithm 1** User-side algorithm of ALLOMFREE.

---
1: **Input :** $\mathbf{v} = [v_1, v_2, ..., v_d]$, $\mathbf{k} = [k_1, k_2, ..., k_d]$, $\mathbb{A} = \{L\text{-}GRR, L\text{-}OSUE\}$,
   $\epsilon_\infty$, $\epsilon_1$, number of reports $\tau$.
2: $r \leftarrow Uniform(\{1, 2, ..., d\})$        ▷ Select attribute only once
3: $B \leftarrow v_r$
4: **if** $Var^*[\hat{f}_{L_{(L\text{-}GRR)}}] \leq Var^*[\hat{f}_{L_{(L\text{-}OSUE)}}]$ :    ▷ Check variances with Eq. (11)
5:      $\mathcal{A} \leftarrow$ L-GRR          ▷ Select L-GRR as local randomizer
6: **else**
7:      $\mathcal{A} \leftarrow$ L-OSUE       ▷ Select L-OSUE as local randomizer
8: $B' \leftarrow \mathcal{A}(B, \epsilon_\infty, k_r)$    ▷ First round of sanitization (permanent memoization)
9: **for** $t \in [1, \tau]$ **do**
10:    $B'' = \mathcal{A}(B', \epsilon_1, k_r)$        ▷ Second round of sanitization
11: **end for**
12: **send :** $(t, \langle r, B'' \rangle)$ for $t \in [1, \tau]$

---

and Multidimensional FREquency Estimates (ALLOMFREE), which is summarized in Algorithm 1 as a pseudocode. The intuition of ALLOMFREE is as follows. By requiring each user to submit only 1 attribute with the whole privacy budget, it reduces the variance incurred (analysis on Section III). Also, since there is a way to measure the approximate variance of the proposed longitudinal protocols (L-GRR and L-OSUE), ALLOMFREE adaptively chooses the protocol with a smaller variance value to optimize the data utility.

According to the analysis in Subsections IV-B and IV-C, Alg. 1 satisfies $\epsilon$-LDP with upper $\epsilon_\infty$ (infinity reports) and lower $\epsilon_1$ (a single report) bounds as it uses either L-GRR or L-OSUE to sanitize a single attribute per user. **Lastly, notice**

**that, to ensure users' privacy over time and to avoid the sequential composition theorem [3], each user must always report the same unique attribute** $A_r$. On the server-side, for each attribute $j \in [1, d]$ at time $t \in [1, \tau]$, the estimated frequency $\hat{f}_L(v_i)$ that a value $v_i$ occurs for $i \in [1, k_j]$ is calculated using Eq. (9).

## V. RESULTS AND DISCUSSION

In this section, we present the setup of our experiments in Subsection V-A, the results with real-world data in Subsection V-B, and a general discussion in Subsection V-C with related work and limitations.

### A. Setup of experiments

The main goal of our experiments is to evaluate the proposed longitudinal LDP protocols on multidimensional frequency estimates a single time, i.e., satisfying $\epsilon_1$-LDP (as in [11], [37], [38], for example).

**Environment.** All algorithms were implemented in Python 3.8.8 with NumPy 1.19.5 and Numba 0.53.1 libraries. The codes we developed and used for all experiments are available in a Github repository[1]. In all experiments, we report average results over 100 runs as LDP algorithms are randomized.

**Methods evaluated.** We consider for evaluation the following solutions and protocols:

- Solution *Smp* (cf. Section III), which randomly samples a single attribute to send with the whole privacy budget. We will experiment with the state-of-the-art protocols, namely, L-SUE and L-OUE, and with our extended protocols L-OSUE and L-SOUE;
- Our ALLOMFREE solution (cf. Alg. 1), which also randomly samples a single attribute to send with the whole privacy budget but adaptively select the optimal protocol, i.e., either L-GRR or L-OSUE.

**Experimental evaluation and metrics.** We vary the longitudinal privacy parameter in the range $\epsilon_\infty = [0.5, 1, ..., 3.5, 4]$ with $\epsilon_1 = [0.3\epsilon_\infty, 0.6\epsilon_\infty]$ to compare our experimental results with numerical ones from Subsection IV-D. Notice that this range of privacy guarantees is commonly used in the literature for multidimensional data (e.g., in [31] the range is $\epsilon = [0.5, ..., 4]$ and in [33] the range is $\epsilon = [0.1, ..., 10]$).

To evaluate our results, we use the mean squared error (MSE) metric averaged per the number of attributes $d$ **in a single data collection** $\tau = 1$, **i.e., with** $\epsilon_1$-**LDP.** Thus, for each attribute $j$, we compute for each value $v(i) \in \mathcal{D}_j$ the estimated frequency $\hat{f}(v_i)$ and the real one $f(v_i)$ and calculate their differences. More precisely,

$$MSE_{avg} = \frac{1}{\tau} \sum_{t \in [1, \tau]} \frac{1}{d} \sum_{j \in [1, d]} \frac{1}{|\mathcal{D}_j|} \sum_{v \in \mathcal{D}_j} (f(v_i) - \hat{f}(v_i))^2.$$
(16)

**Datasets.** For ease of reproducibility, we conduct our experiments on four multidimensional open datasets.

---

- *Nursery.* A dataset from the UCI machine learning repository [43] with $d = 9$ categorical attributes and $n = 12960$ samples. The domain size of each attribute is $\mathbf{k} = [3, 5, 4, 4, 3, 2, 3, 3, 5]$, respectively.
- *Adult.* A dataset from the UCI machine learning repository [43] with $d = 9$ categorical attributes and $n = 45222$ samples after cleaning the data. The domain size of each attribute is $\mathbf{k} = [7, 16, 7, 14, 6, 5, 2, 41, 2]$, respectively.
- *MS-FIMU.* An open dataset from [44] with $d = 6$ categorical attributes and $n = 88935$ samples. The domain size of each attribute is $\mathbf{k} = [3, 3, 8, 12, 37, 11]$, respectively.
- *Census-Income.* A dataset from the UCI machine learning repository [43] with $d = 33$ categorical attributes and $n = 299285$ samples. The domain size of each attribute is $\mathbf{k} = [9, 52, 47, 17, 3, ..., 43, 43, 43, 5, 3, 3, 3, 2]$, respectively.

### B. Results

Our experiments were conducted on four real-world datasets with varied parameters for $n$, $d$, and $\mathbf{k}$, which allowed evaluating our solutions more practically. Fig. 4 (*Nursery*), Fig. 5 (*Adult*), Fig. 6 (*MS-FIMU*), and Fig. 7 (*Census-Income*) illustrate for all evaluated protocols, averaged $MSE_{avg}$ (y-axis) according to the longitudinal privacy parameter $\epsilon_\infty$ (x-axis) with $\epsilon_1 = 0.3\epsilon_\infty$ (plot on the top) and with $\epsilon_1 = 0.6\epsilon_\infty$ (plot on the bottom), respectively.

As one can notice in the results, for all datasets, ALLOM-FREE consistently and considerably outperforms the state-of-the-art protocols, namely, L-SUE (a.k.a. Basic-RAPPOR) [11] and L-OUE (that uses OUE [14] twice). Indeed, the difference on performance between ALLOMFREE and the other longitudinal LDP protocols increases proportionally according to the privacy guarantees, i.e., for high $\epsilon_\infty$ and $\epsilon_1$ values the gap is bigger. This is, first, because in all datasets there are attribute(s) with small domain size (e.g., $k_j = 2$ or $k_j = 3$), in which L-GRR can provide smaller variance values than L-UE protocols (cf. Subsection IV-D). Secondly, by selecting adequately the probabilities $p_1, q_1, p_2, q_2$ for the L-UE protocol (i.e., L-OSUE) also optimizes data utility. Thus, since there is a way to measure the approximate variance of the extended protocols (i.e., Eq. (11)), given the sampled attribute, ALLOM-FREE adaptively selects one of the optimized protocol (i.e., L-GRR or L-OSUE) whose smaller variance improves the data utility.

In addition, among the L-UE protocols applied individually, the experimental results with multidimensional data approximate the numerical results with a single attribute from Subsection IV-D. For instance, the proposed L-OSUE provides similar or better performance than L-SUE while always outperforming L-OUE. Besides, L-SOUE always outperforms L-OUE too, achieving similar performance than L-OSUE and L-SUE in low privacy regimes (i.e., high $\epsilon$ values). As we have already shown in Subsection IV-D, even though OUE has better utility than SUE for one-time collection [14], applying OUE twice does not provide higher utility.
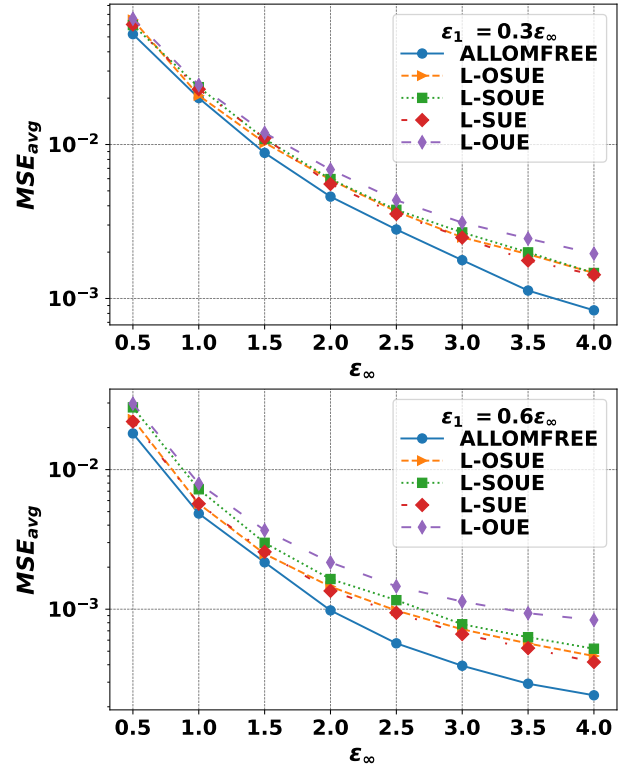


Fig. 4: Averaged MSE varying $\epsilon_\infty$ with $\epsilon_1 = 0.3\epsilon_\infty$ (plot on the top) and with $\epsilon_1 = 0.6\epsilon_\infty$ (plot on the bottom) on the *Nursery* dataset.
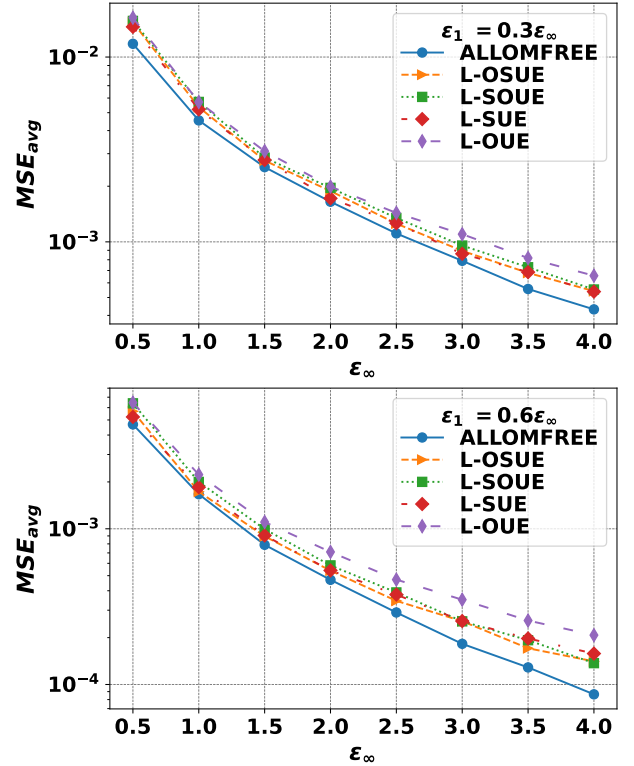


Fig. 5: Averaged MSE varying $\epsilon_\infty$ with $\epsilon_1 = 0.3\epsilon_\infty$ (plot on the top) and with $\epsilon_1 = 0.6\epsilon_\infty$ (plot on the bottom) on the *Adult* dataset.

Fig. 6: Averaged MSE varying $\epsilon_\infty$ with $\epsilon_1 = 0.3\epsilon_\infty$ (plot on the top) and with $\epsilon_1 = 0.6\epsilon_\infty$ (plot on the bottom) on the *MS-FIMU* dataset.


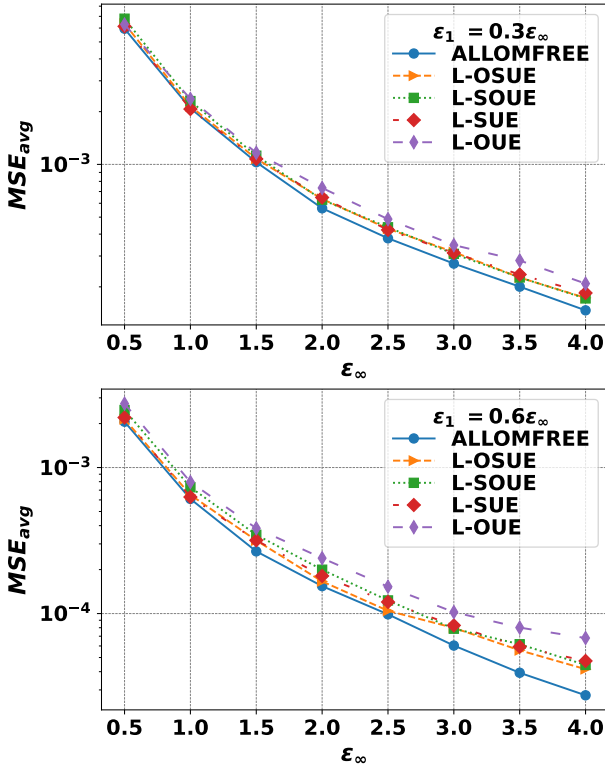
Fig. 7: Averaged MSE varying $\epsilon_\infty$ with $\epsilon_1 = 0.3\epsilon_\infty$ (plot on the top) and with $\epsilon_1 = 0.6\epsilon_\infty$ (plot on the bottom) on the *Census-Income* dataset.
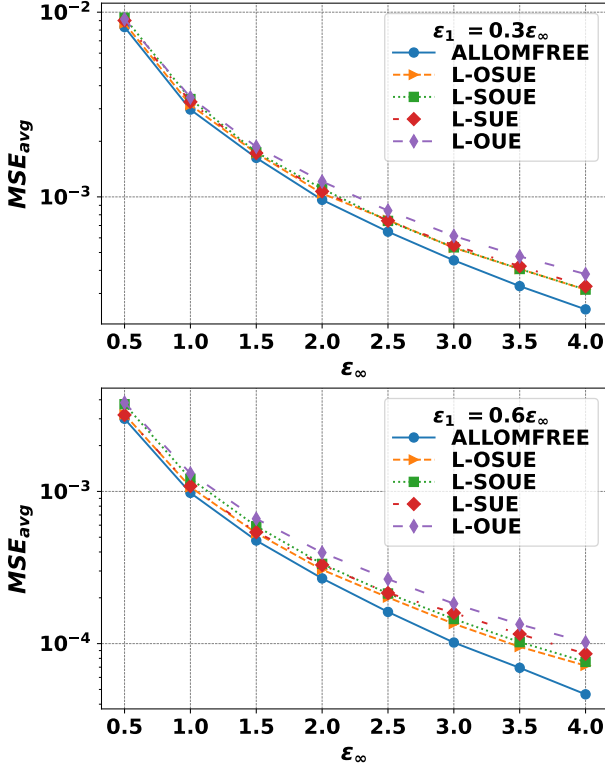
| $\epsilon_\infty$ | Nursery | | Adult | | MS-FIMU | | Census-Income | |
|---|---|---|---|---|---|---|---|---|
| | $\mathcal{U}_{L\text{-}SUE}$ | $\mathcal{U}_{L\text{-}OUE}$ | $\mathcal{U}_{L\text{-}SUE}$ | $\mathcal{U}_{L\text{-}OUE}$ | $\mathcal{U}_{L\text{-}SUE}$ | $\mathcal{U}_{L\text{-}OUE}$ | $\mathcal{U}_{L\text{-}SUE}$ | $\mathcal{U}_{L\text{-}OUE}$ |
| 0.5 | 13.51 | 20.63 | 19.03 | 27.73 | 3.03 | 5.43 | 7.84 | 9.48 |
| 1.0 | 12.36 | 17.75 | 12.77 | 20.44 | 1.01 | 11.57 | 9.21 | 14.08 |
| 1.5 | 19.95 | 25.86 | 8.47 | 18.01 | 4.13 | 11.55 | 5.82 | 12.92 |
| 2.0 | 17.18 | 33.24 | 4.11 | 17.16 | 13.22 | 23.44 | 10.06 | 20.41 |
| 2.5 | 20.70 | 35.40 | 11.93 | 22.54 | 10.41 | 22.25 | 12.77 | 23.15 |
| 3.0 | 28.69 | 42.98 | 8.35 | 28.22 | 13.07 | 21.56 | 17.07 | 26.21 |
| 3.5 | 36.19 | 54.02 | 18.97 | 32.02 | 14.78 | 29.10 | 22.02 | 30.96 |
| 4.0 | 41.24 | 57.16 | 19.81 | 34.25 | 20.38 | 29.64 | 24.99 | 35.60 |
| Mean | 23.73 | 35.88 | 12.93 | 25.05 | 10.00 | 19.32 | 13.72 | 21.60 |

TABLE III: Accuracy gain of ALLOMFREE over the state-of-the-art L-SUE and L-OUE protocols for all datasets with $\epsilon_1 = 0.3\epsilon_\infty$, measured with the $\mathcal{U}_{L\text{-}SUE}$ and $\mathcal{U}_{L\text{-}OUE}$ metrics expressed in %.

| $\epsilon_\infty$ | Nursery | | Adult | | MS-FIMU | | Census-Income | |
|---|---|---|---|---|---|---|---|---|
| | $\mathcal{U}_{L\text{-}SUE}$ | $\mathcal{U}_{L\text{-}OUE}$ | $\mathcal{U}_{L\text{-}SUE}$ | $\mathcal{U}_{L\text{-}OUE}$ | $\mathcal{U}_{L\text{-}SUE}$ | $\mathcal{U}_{L\text{-}OUE}$ | $\mathcal{U}_{L\text{-}SUE}$ | $\mathcal{U}_{L\text{-}OUE}$ |
| 0.5 | 17.82 | 38.84 | 10.42 | 27.46 | 6.41 | 24.79 | 5.65 | 21.61 |
| 1.0 | 14.99 | 38.97 | 9.83 | 25.14 | 2.97 | 23.32 | 9.79 | 25.46 |
| 1.5 | 15.88 | 41.05 | 12.90 | 28.59 | 16.00 | 30.52 | 11.88 | 28.05 |
| 2.0 | 27.52 | 54.69 | 12.95 | 33.78 | 14.81 | 35.65 | 18.45 | 32.31 |
| 2.5 | 39.59 | 60.96 | 23.28 | 38.50 | 17.71 | 35.34 | 24.89 | 39.11 |
| 3.0 | 40.64 | 65.32 | 28.59 | 47.95 | 27.26 | 40.97 | 36.12 | 44.48 |
| 3.5 | 44.39 | 68.73 | 34.85 | 50.00 | 33.69 | 50.94 | 40.01 | 48.18 |
| 4.0 | 42.24 | 71.13 | 45.26 | 58.33 | 41.83 | 59.47 | 45.85 | 54.44 |
| Mean | 30.38 | 54.96 | 22.26 | 38.72 | 20.08 | 37.62 | 24.08 | 36.70 |

TABLE IV: Accuracy gain of ALLOMFREE over the state-of-the-art L-SUE and L-OUE protocols for all datasets with $\epsilon_1 = 0.6\epsilon_\infty$, measured with the $\mathcal{U}_{L\text{-}SUE}$ and $\mathcal{U}_{L\text{-}OUE}$ metrics expressed in %.

To complement the results of Figs. 4 – 7, Table III ($\epsilon_1 = 0.3\epsilon_\infty$) and Table IV ($\epsilon_1 = 0.6\epsilon_\infty$) exhibit for all datasets and $\epsilon_\infty$ guarantees the following utility metrics:

$$\mathcal{U}_{L\text{-}SUE} = \frac{MSE_{avg(L\text{-}SUE)} - MSE_{avg(ALLOMFREE)}}{MSE_{avg(L\text{-}SUE)}},$$
$$\mathcal{U}_{L\text{-}OUE} = \frac{MSE_{avg(L\text{-}OUE)} - MSE_{avg(ALLOMFREE)}}{MSE_{avg(L\text{-}OUE)}}, \quad (17)$$

in which $\mathcal{U}_{L\text{-}SUE}$ and $\mathcal{U}_{L\text{-}OUE}$ represent the accuracy gain of ALLOMFREE over the state-of-the-art L-SUE and L-OUE protocols, respectively.

From Tables III and IV, one can notice that ALLOMFREE considerably improves the quality of the frequency estimates in comparison with the state-of-the-art L-SUE and L-OUE protocols. On average, ALLOMFREE improves the results of L-SUE at least 10% with the *MS-FIMU* dataset in Table III and at most 30.38% with the *Nursery* dataset in Table IV for the privacy guarantees $\epsilon_\infty$ and $\epsilon_1$ analyzed. Similarly, on average, ALLOMFREE improves the results of L-OUE at least 19.32% with the *MS-FIMU* dataset in Table III and at most 54.96% with the *Nursery* dataset in Table IV. The highest gain of accuracy was about $\sim 71\%$, achieved with the *Nursery*

dataset when $\epsilon_\infty = 4$ in Table IV in comparison with the L-OUE protocol. Finally, as one can note, with higher values of $\epsilon_1$, ALLOMFREE will provide much higher utility than the other protocols.

### C. Discussion

In recent times, there have been several works on the local DP setting in both academia [10], [14], [16]–[18], [30], [31], [33], [45], [46] and practical deployment [11]–[13], [47]. The local DP model does not rely on collecting raw data anymore, which has a clear connection with the concept of randomized response [39]. Among many other complex tasks (e.g., heavy hitter estimation [35], [42], [45], marginal estimation [25]–[29], analytical/range queries [21]–[24], frequent itemset mining [40], [48]), frequency estimation is a fundamental primitive in LDP and has received considerable attention for a single attribute [11], [12], [14]–[20], [33], [37], [46], [49].

However, most studies for collecting multidimensional data with LDP mainly focused on numerical data [46] (e.g., [30]–[33]) or other complex tasks with categorical data (e.g., marginal estimation [25]–[29], analytical/range queries [21]–[24]). Besides, most frequency estimation academic literature focuses on single data collection. To address longitudinal data collections, in [11], [12], the authors proposed LDP protocols based on two rounds of sanitization, i.e., *memoization*, which was also adopted in this paper. In the literature, some works [37], [38] applied L-SUE (a.k.a. Basic-RAPPOR [11]) and L-OUE (i.e., OUE [14] with memoization) for longitudinal frequency estimates. However, rather than strictly using only SUE or OUE, we prove that the optimal combination is starting with OUE and then with SUE (i.e., L-OSUE). Indeed, both "multiple" settings combined (i.e., many attributes and several collections throughout time), imposes several challenges, in which this paper, proposes the first solution named ALLOMFREE under LDP.

Lastly, some limitations and prospective directions of this paper are described in the following. Concerning the privacy guarantees, the memoization step of ALLOMFREE is certainly effective for longitudinal privacy to the cases where the true client's data does not vary (static) or vary very slowly or in an uncorrelated manner [11]. In many application scenarios, gender, age-ranges, nationality, and other demographic data are generally static or vary hardly ever. On the other hand, for dynamic attributes such as location or the time spent in the application, this is not the case. Therefore, for each different value, a new memoized value would be generated, thus accumulating the privacy budget $\epsilon_\infty$ by the sequential composition theorem [3]. Besides, ALLOMFREE is based on the multidimensional *Smp* solution, which randomly samples a single attribute per user out of $d$ ones. However, aggregators (who are also seen as attackers) are aware of the sampled attribute and its LDP value, which is protected by a "less strict" $e^\epsilon$ probability bound (rather than $e^{\epsilon/d}$). Indeed, in some cases, using the *Smp* solution (as well as ALLOMFREE) may be unfair with some users, e.g., users that randomly sample a demographic attribute (e.g., age) might be less concerned

to report their data than those whose sampled attribute is socially "more" sensitive (e.g., disease or location). Investigating how to deal with these issues is a prospective direction. Thus, besides the aforementioned directions, for future work, we also suggest and intend to improve frequency estimates through post-processing techniques [41], [50] and to design LDP protocols for longitudinal and multidimensional studies considering both numerical and categorical data.

## VI. CONCLUSION

This paper investigates the problem of collecting multi-dimensional data throughout time for the fundamental task of frequency estimation under LDP guarantees. We extended and analyzed three state-of-the-art LDP protocols, namely, GRR [16], OUE [14], and SUE [11], and proposed an optimized solution namely ALLOMFREE, which randomly samples one attribute per user and adaptively selects a protocol with lower variance (i.e., L-GRR or L-OSUE) in order to improve data utility. We demonstrate through experimental validations using four real-world datasets the advantages of ALLOMFREE over the state-of-the-art protocols L-SUE [11] and L-OUE [14], with gain of accuracy, on average, ranging from $10\%$ up to $55\%$ with the analyzed range of $\epsilon_\infty$ and $\epsilon_1$ privacy guarantees.

## REFERENCES

[1] C. Dwork, F. McSherry, K. Nissim, and A. Smith, "Calibrating noise to sensitivity in private data analysis," in *Theory of Cryptography*. Springer Berlin Heidelberg, 2006, pp. 265–284.

[2] C. Dwork, "Differential privacy," in *Automata, Languages and Programming*. Springer Berlin Heidelberg, 2006, pp. 1–12. [Online]. Available: https://doi.org/10.1007/11787006_1

[3] C. Dwork, A. Roth *et al.*, "The algorithmic foundations of differential privacy," *Foundations and Trends® in Theoretical Computer Science*, vol. 9, no. 3–4, pp. 211–407, 2014.

[4] A. Aktay *et al.*, "Google COVID-19 community mobility reports: Anonymization process description (version 1.0)," *arXiv preprint arXiv:2004.04145*, 2020.

[5] R. Rogers, S. Subramaniam, S. Peng, D. Durfee, S. Lee, S. K. Kancha, S. Sahay, and P. Ahammad, "Linkedin's audience engagements API: A privacy preserving data analytics system at scale," *arXiv preprint arXiv:2002.05839*, 2020.

[6] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, "Deep learning with differential privacy," ser. CCS '16. New York, NY, USA: Association for Computing Machinery, 2016, p. 308–318.

[7] J. M. Abowd, "The U.S. census bureau adopts differential privacy," in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, Jul. 2018.

[8] S. Garfinkel, "Implementing differential privacy for the 2020 census." USENIX Association, Feb. 2021.

[9] D. McCandless, T. Evans, M. Quick, E. Hollowood, C. Miles, D. Hampson, and D. Geere, "World's biggest data breaches & hacks," https://www.informationisbeautiful.net/visualizations/worlds-biggest-data-breaches-hacks/, jan 2021, online; accessed 11 March 2021.

[10] S. P. Kasiviswanathan, H. K. Lee, K. Nissim, S. Raskhodnikova, and A. Smith, "What can we learn privately?" in *2008 49th Annual IEEE Symposium on Foundations of Computer Science.* IEEE, Oct. 2008.

[11] U. Erlingsson, V. Pihur, and A. Korolova, "RAPPOR: Randomized aggregatable privacy-preserving ordinal response," in *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security.* New York, NY, USA: ACM, 2014, pp. 1054–1067.

[12] B. Ding, J. Kulkarni, and S. Yekhanin, "Collecting telemetry data privately," in *Advances in Neural Information Processing Systems 30*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds. Curran Associates, Inc., 2017, pp. 3571–3580.

[13] A. D. P. Team, "Learning with privacy at scale," https://docs-assets.developer.apple.com/ml-research/papers/learning-with-privacy-at-scale.pdf, dec 2017, online; accessed 11 March 2021.

[14] T. Wang, J. Blocki, N. Li, and S. Jha, "Locally differentially private protocols for frequency estimation," in *26th USENIX Security Symposium (USENIX Security 17).* Vancouver, BC: USENIX Association, Aug. 2017, pp. 729–745.

[15] S. Wang, Y. Nie, P. Wang, H. Xu, W. Yang, and L. Huang, "Local private ordinal data distribution estimation," in *IEEE INFOCOM 2017 - IEEE Conference on Computer Communications.* IEEE, May 2017.

[16] P. Kairouz, K. Bonawitz, and D. Ramage, "Discrete distribution estimation under local privacy," in *International Conference on Machine Learning.* PMLR, 2016, pp. 2436–2444.

[17] J. Acharya, Z. Sun, and H. Zhang, "Hadamard response: Estimating distributions privately, efficiently, and with little communication," in *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, ser. Proceedings of Machine Learning Research, K. Chaudhuri and M. Sugiyama, Eds., vol. 89. PMLR, 16–18 Apr 2019, pp. 1120–1129.

[18] M. Alvim, K. Chatzikokolakis, C. Palamidessi, and A. Pazii, "Invited paper: Local differential privacy on metric spaces: Optimizing the trade-off with utility," in *2018 IEEE 31st Computer Security Foundations Symposium (CSF).* IEEE, Jul. 2018.

[19] D. Zhao, H. Chen, S. Zhao, X. Zhang, C. Li, and R. Liu, "Local differential privacy with k-anonymous for frequency estimation," in *2019 IEEE International Conference on Big Data (Big Data).* IEEE, Dec. 2019.

[20] Z. Li, T. Wang, M. Lopuhaä-Zwakenberg, N. Li, and B. Škoric, "Estimating numerical distributions under local differential privacy," in *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data.* ACM, May 2020.

[21] M. Xu, B. Ding, T. Wang, and J. Zhou, "Collecting and analyzing data jointly from multiple services under local differential privacy," *Proceedings of the VLDB Endowment*, vol. 13, no. 12, pp. 2760–2772, Aug. 2020.

[22] J. Yang, T. Wang, N. Li, X. Cheng, and S. Su, "Answering multi-dimensional range queries under local differential privacy," *Proc. VLDB Endow.*, vol. 14, no. 3, p. 378–390, Nov. 2020.

[23] X. Gu, M. Li, Y. Cao, and L. Xiong, "Supporting both range queries and frequency estimation with local differential privacy," in *2019 IEEE Conference on Communications and Network Security (CNS).* IEEE, Jun. 2019.

[24] G. Cormode, T. Kulkarni, and D. Srivastava, "Answering range queries under local differential privacy," *Proceedings of the VLDB Endowment*, vol. 12, no. 10, pp. 1126–1138, Jun. 2019.

[25] Z. Shen, Z. Xia, and P. Yu, "PLDP: Personalized local differential privacy for multidimensional data aggregation," *Security and Communication Networks*, vol. 2021, pp. 1–13, Jan. 2021.

[26] F. Peng, S. Tang, B. Zhao, and Y. Liu, "A privacy-preserving data aggregation of mobile crowdsensing based on local differential privacy," in *Proceedings of the ACM Turing Celebration Conference - China.* ACM, May 2019.

[27] Z. Zhang, T. Wang, N. Li, S. He, and J. Chen, "CALM: Consistent adaptive local marginal for marginal release under local differential privacy," *Proceedings of the ACM Conference on Computer and Communications Security*, pp. 212–229, 2018.

[28] X. Ren, C.-m. Yu, W. Yu, S. Yang, S. Member, X. Yang, J. A. Mccann, P. S. Yu, and L. Fellow, "LoPub : High-Dimensional Crowdsourced Data," vol. 13, no. 9, pp. 2151–2166, 2018.

[29] G. Fanti, V. Pihur, and Ú. Erlingsson, "Building a RAPPOR with the unknown: Privacy-preserving learning of associations and data dictionaries," *Proceedings on Privacy Enhancing Technologies*, vol. 2016, no. 3, pp. 41–61, May 2016.

[30] T. T. Nguyên, X. Xiao, Y. Yang, S. C. Hui, H. Shin, and J. Shin, "Collecting and analyzing data from smart device users with local differential privacy." *ArXiv*, vol. abs/1606.05053, 2016.

[31] N. Wang, X. Xiao, Y. Yang, J. Zhao, S. C. Hui, H. Shin, J. Shin, and G. Yu, "Collecting and analyzing multidimensional data with local differential privacy," in *2019 IEEE 35th International Conference on Data Engineering (ICDE).* IEEE, Apr. 2019.

[32] J. C. Duchi, M. I. Jordan, and M. J. Wainwright, "Minimax optimal procedures for locally private estimation," *Journal of the American Statistical Association*, vol. 113, no. 521, pp. 182–201, Jan. 2018.

[33] T. Wang, J. Zhao, Z. Hu, X. Yang, X. Ren, and K.-Y. Lam, "Local differential privacy for data collection and analysis," *Neurocomputing*, vol. 426, pp. 114–133, Feb. 2021.

[34] Ú. Erlingsson, V. Feldman, I. Mironov, A. Raghunathan, S. Song, K. Talwar, and A. Thakurta, "Encode, shuffle, analyze privacy revisited: Formalizations and empirical evaluation," *arXiv preprint arXiv:2001.03618*, 2020.

[35] T. Wang, N. Li, and S. Jha, "Locally differentially private heavy hitter identification," *IEEE Transactions on Dependable and Secure Computing*, vol. 18, no. 2, pp. 982–993, Mar. 2021.

[36] H. H. Arcolezi, J.-F. Couchot, B. A. Bouna, and X. Xiao, "Longitudinal collection and analysis of mobile phone data with local differential privacy," in *IFIP Advances in Information and Communication Technology.* Springer International Publishing, 2021, pp. 40–57.

[37] J. W. Kim, D.-H. Kim, and B. Jang, "Application of local differential privacy to collection of indoor positioning data," *IEEE Access*, vol. 6, pp. 4276–4286, 2018.

[38] I. D. C. Vidal, A. L. da Costa Mendonça, F. Rousseau, and J. D. C. Machado, "ProTECting: An application of local differential privacy for IoT at the edge in smart home scenarios," in *Anais XXXVIII Simpósio Brasileiro de Redes de Computadores e Sistemas Distribuídos (SBRC 2020).* Sociedade Brasileira de Computação, Dec. 2020.

[39] S. L. Warner, "Randomized response: A survey technique for eliminating evasive answer bias," *Journal of the American Statistical Association*, vol. 60, no. 309, pp. 63–69, Mar. 1965.

[40] T. Wang, N. Li, and S. Jha, "Locally differentially private frequent itemset mining," in *2018 IEEE Symposium on Security and Privacy (SP).* IEEE, May 2018.

[41] T. Wang, M. Lopuhaa-Zwakenberg, Z. Li, B. Skoric, and N. Li, "Locally differentially private frequency estimation with consistency," in *Proceedings 2020 Network and Distributed System Security Symposium.* Internet Society, 2020.

[42] R. Bassily, K. Nissim, U. Stemmer, and A. Thakurta, "Practical locally private heavy hitters," *arXiv preprint arXiv:1707.04982*, 2017.

[43] D. Dua and C. Graff, "UCI machine learning repository," 2017. [Online]. Available: http://archive.ics.uci.edu/ml

[44] H. H. Arcolezi, J.-F. Couchot, O. Baala, J.-M. Contet, B. A. Bouna, and X. Xiao, "Mobility modeling through mobile data: generating an optimized and open dataset respecting privacy," in *2020 International Wireless Communications and Mobile Computing (IWCMC).* IEEE, Jun. 2020.

[45] R. Bassily and A. Smith, "Local, private, efficient protocols for succinct histograms," in *Proceedings of the forty-seventh annual ACM symposium on Theory of Computing.* ACM, Jun. 2015.

[46] X. Xiong, S. Liu, D. Li, Z. Cai, and X. Niu, "A comprehensive survey on local differential privacy," *Security and Communication Networks*, vol. 2020, pp. 1–29, Oct. 2020.

[47] S. Kessler, J. Hoff, and J.-C. Freytag, "SAP HANA goes private," *Proceedings of the VLDB Endowment*, vol. 12, no. 12, pp. 1998–2009, Aug. 2019.

[48] 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 *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security.* ACM, Oct. 2016.

[49] H. H. Arcolezi, J.-F. Couchot, S. Cerna, C. Guyeux, G. Royer, B. A. Bouna, and X. Xiao, "Forecasting the number of firefighter interventions per region with local-differential-privacy-based data," *Computers & Security*, vol. 96, p. 101888, Sep. 2020.

[50] E. ElSalamouny and C. Palamidessi, "Generalized iterative bayesian update and applications to mechanisms for privacy protection," in *2020 IEEE European Symposium on Security and Privacy (EuroS&P).* IEEE, Sep. 2020.