

Differentially Private Aggregation of Distributed Time-Series with Transformation and Encryption

Vibhor Rastogi *
University of Washington
Seattle, WA, USA, 98105-2350
vibhor@cs.washington.edu

Suman Nath
Microsoft Research
Redmond, WA, USA, 98052-6399
sumann@microsoft.com

ABSTRACT

We propose PASTE, the first differentially private aggregation algorithms for distributed time-series data that offer good practical utility without any trusted server. PASTE addresses two important challenges in participatory data-mining applications where (i) individual users collect temporally correlated time-series data (such as location traces, web history, personal health data), and (ii) an untrusted third-party aggregator wishes to run aggregate queries on the data. To address this, PASTE incorporates two new algorithms.

To ensure differential privacy for time-series data despite the presence of temporal correlation, PASTE uses the Fourier Perturbation Algorithm (FPA_k). Standard differential privacy techniques perform poorly for time-series data. To answer n queries, such techniques can result in a noise of $\Theta(n)$ to each query answer, making the answers practically useless if n is large. Our FPA_k algorithm perturbs the Discrete Fourier Transform of the query answers. For answering n queries, FPA_k improves the expected error from $\Theta(n)$ to roughly $\Theta(k)$ where k is the number of Fourier coefficients that can (approximately) reconstruct all the n query answers. Our experiments show that $k \ll n$ for many real-life-data-sets resulting in a huge error-improvement for FPA_k.

To deal with the absence of a trusted central server, PASTE uses the Distributed Laplace Perturbation Algorithm (DLPA) that adds noise in a distributed way in order to guarantee differential privacy. To the best of our knowledge, DLPA is the first distributed differentially private algorithm that can scale with a large number of users: DLPA outperforms the only other distributed solution for differential privacy proposed so far, by reducing the computational load per user from $O(U)$ to $O(1)$ where U is the number of users.

Categories and Subject Descriptors

H.2.8 [DATABASE MANAGEMENT]: Database applications—*Data Mining*; C.2.4 [DISTRIBUTED SYSTEMS]: Distributed applications; G.3 [PROBABILITY AND STATISTICS]: Time series analysis

*Work done while visiting Microsoft Research

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SIGMOD'10, June 6–11, 2010, Indianapolis, Indiana, USA.
Copyright 2010 ACM 978-1-4503-0032-2/10/06 ...\$10.00.

General Terms

Algorithms, Performance, Security

Keywords

private data analysis, participatory data mining, output perturbation, time-series data, distributed noise addition

1. INTRODUCTION

The ever increasing instrumentation of the physical and the virtual world has given us an unprecedented opportunity to collect useful data from diverse sources and to mine it for understanding important phenomena. Consider the following examples of *participatory data mining applications*.

E1: In participatory sensing applications such as CarTel [16], BikeNet [8], PEIR [25], WeatherUnderground¹, participants contribute various time-series data, e.g., their current locations, speeds, weather information, images, etc. These data can be aggregated and mined for useful information such as community interests (e.g., popular places), congestion patterns in roads, micro-weather, etc.

E2: A Web browser can install plug-ins to monitor users' browsing behaviors such as the numbers of different types of web pages a user visits, the types of products he buys from online stores, etc. Historical data from such plug-ins can be aggregated to understand user behaviors for improving search results or for better targeted advertisement delivery [14].

E3: Health-care applications such as Microsoft HealthVault² and Google Health³ allow users to periodically upload data such as their weights, occurrences of diseases, amounts of exercise, food and drug intake, etc. PIER [25] allows individual users to store such data locally in *personal data vaults*. Such data can be mined in combination to understand disease outbreaks, distribution of weights, relationship of weight gain and drug intake, etc.

As the above examples suggest, aggregate statistics computed from data contributed by a large number of individual participants can be quite useful. However, data owners or publishers may not be always willing to reveal the true values due to various reasons, most notably privacy considerations. The goal of our work is to enable third parties to compute useful aggregate queries over, while guaranteeing the privacy of, the data from individual publishers.

Random perturbation is a widely employed and accepted approach for partial information hiding, which works by introducing uncertainty about individual values. Prior perturbation techniques are designed primarily for relational data [7, 9, 19, 26]. However,

¹<http://www.weatherunderground.com>

²<http://www.healthvault.com>

³<https://www.google.com/health>

	Centralized	Distributed
Relational data	E.g., [7, 15, 27] <i>Inaccurate for long time-series query sequences</i>	E.g., [6, 9, 12, 28] <i>Either inefficient or inaccurate [9] for a large number of users</i>
Time-series data	E.g., [11, 20] <i>No formal privacy guarantees</i>	PASTE (this paper) <i>Accurate & efficient answers for time-series query sequences under differential privacy</i>

Table 1: The design space and existing works

participatory data mining applications have two unique characteristics that make existing privacy mechanisms inadequate for these applications.

- **Time-series data:** The applications generate time series numerical or categorical data. Data at successive timestamps from the same source can be highly correlated.
- **Distributed sources:** Data publishers may not trust any single third party to see their true data. This means, the querier needs to be able to compute useful aggregates without seeing the true data values.

The above characteristics make most existing privacy solutions, which assume relational data with negligible correlations across tuples [9, 19, 26] or existence of a central trusted entity for carefully introducing noise [7, 15, 27], inadequate for our target applications (as summarized in Table 1).

Thus, to realize widespread adoption of participatory data mining applications, one needs to address two challenges. The first challenge is to ensure privacy for time-series data, which is problematic due to the strong correlation among successive values in the series. This correlation makes answers to different queries over time-series data to also become correlated, e.g. a sequence of queries computing the average weight of a community at successive weeks.

One possible way for achieving privacy is to perturb the answers to such queries independently of one another, thereby ensuring that even revealing a few true answers does not help infer anything about the perturbation of other answers. However, [20] pointed out that if the time-series exhibit certain patterns, then independent perturbation of query answers can be distinguished from the original answers and filtered out. Authors in [11, 20] consider perturbing time series data to defend against several privacy attacks, but they do not provide any formal privacy guarantee, without which data owners may not publish sensitive data in the fear of unforeseen privacy attacks.

On the other hand, formal privacy guarantees like differential privacy that work well for relational data, seem too hard to achieve for time series data. For instance, standard differentially private techniques [7] can result in a noise of $\Theta(n)$ to each query answer, where n is the number of queries to answer, making the query answers practically useless if a long sequence of queries is to be answered.

The second challenge arises from the absence of a trusted aggregator. Most previous works assume that a trusted aggregator, who has access to the raw data, computes target functions on the data and then perturbs the results [7, 19, 26]. In the absence of a trusted aggregator, users need to perturb their data before publishing it to the aggregator [9]. However, if users perturb data independently, the noise variance in the perturbed estimate grows linearly with the number of users, reducing the utility of the aggregate information. To improve utility, cryptographic techniques like Secure Multiparty

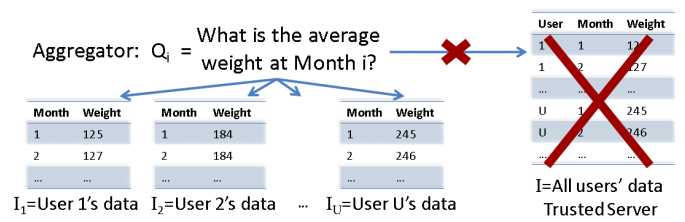


Figure 1: System Model (Users with data I_1, \dots, I_U Aggregator issues recurring query $Q = Q_1, \dots, Q_n$ No trusted server has $I = I_1 \cup I_2 \dots \cup I_U$ to evaluate $Q(I)$)

Computation [6, 12, 28] can be used to compute accurate perturbed estimates in a distributed setting. However the computational performance of such cryptographic techniques does not scale well with a large number of users.

In this paper, we propose PASTE (Private Aggregation of Signals with Transformation and Encryption), a suite of algorithms that address these two challenges of distributed time-series data. We use the state-of-the-art differential privacy as the privacy requirement, and make the following contributions:

- To answer multiple queries over time-series data under differential privacy, we propose the FPA $_k$ algorithm that perturbs the Discrete Fourier Transform (DFT) of the query answers. For answering n queries, FPA $_k$ improves the error from $\Theta(n)$ (error of standard differential privacy techniques) to roughly $\Theta(k)$ where k is the number of DFT coefficients that can (approximately) reconstruct all the n query answers. Our experiments show that a small $k \ll n$ is sufficient for many real-life datasets, resulting in a huge error-improvement for FPA $_k$. To the best of our knowledge, FPA $_k$ is the first differentially private technique (unlike [11, 20]) that offers practical utility for time-series data.

- We propose the DLPA algorithm that adds noise in a distributed way for providing differential privacy. To the best of our knowledge, DLPA is the first distributed differentially private algorithm that scales with a large number of users: DLPA outperforms the only other proposed distributed algorithm [6], by reducing the computational load per user from $O(U)$ to $O(1)$ where U is the number of users.

- PASTE, our distributed solution, combines the FPA $_k$ and DLPA algorithms to get the accuracy benefits of the former and the scalability of the latter. We empirically evaluate our solution over three real time-series datasets, namely, *GPS traces*, *daily body-weight readings*, and *traffic volumes*. Our experiments show that our solution improves accuracy of query answers by orders of magnitude and also scales well with a large number of users.

We believe that PASTE is an important first step towards practical participatory data mining applications. We have implemented some of our techniques in Microsoft SensorMap,⁴ a real online participatory sensing application that has been publicly available for last three years with several hundreds data publishers. PASTE now allow users to publish private data without revealing the true values and SensorMap to compute useful aggregates over private data.

2. PRELIMINARIES

2.1 Problem Setup and System Model

Motivated by the participatory applications in Section 1, we consider a system model as shown in Figure 1. The system has two types of parties involved: a set of U users and an aggregator. The

⁴<http://www.sensormap.org>

figure shows U users locally storing their personal weight time-series data. We will use the weight time-series as a running example throughout the paper. In general, we model each user u 's data as a (uni- or multi-variate) time series data, and denote it as I_u . We also denote $I = I_1 \cup I_2 \dots \cup I_U$ the combined time-series data of all users. There is no trusted central server, and hence I is never computed. The aggregator, however, wishes to compute aggregate queries over I .

Types of queries. An aggregate query can be a *snapshot* query that returns a single number, e.g., the current average weight over all users. A query can also be *recurring* and ask for periodic answers, e.g., the average weight of all users computed once every month of the year. We model a recurring query as a sequence of queries $\mathbf{Q} = \{Q_1, \dots, Q_n\}$ where each Q_i is a snapshot query. We denote $Q_i(I)$ the value of the snapshot query Q_i on input I , and $\mathbf{Q}(I)$ the vector $\{Q_1(I), \dots, Q_n(I)\}$. A recurring query can be *historical* that focuses on past data only, or *real-time* that runs on data as it arrives.

The data-mining applications we consider are primarily concerned with historical recurring queries, and we develop new algorithms to answer them accurately. Note that even though we consider time-series data in the paper, our techniques can be applied to any types of input data I as long as the query answer $q_i(I)$ is numerical and hence the vector $Q(I)$ is a sequence of numeric values.

Distributed Computation. If a query Q contains predicates on public static attributes, it can be forwarded only to users satisfying the predicates. Otherwise, Q needs to be forwarded to all users. Upon receiving the query Q , a user evaluates Q on his own time series, I_u , perturbs the result, and sends the perturbed results back to the aggregator. The aggregator combines the perturbed results from all users to produce the final result.

A prerequisite for such a distributed system is that the true query answer, $Q(I)$, is computable distributedly. Of all such queries, we consider only queries of the general form $Q(I) = \sum_u f_u(I_u)$, where f_u is an arbitrary function that maps user u 's data, I_u , to numbers.⁵ Such queries, called aggregate-sum queries, are quite general, and as explained in [3], are powerful tools for learning and statistical analysis: many algorithms like correlation, PCA, SVD, decision tree classification, etc. can be implemented using only aggregate-sum queries as a primitive. Queries not included in this class are queries that require a non-sum function to be evaluated collectively over multiple users' data (e.g., aggregate-max or min queries).

Attack Model. We allow both users and the aggregator to be malicious. A malicious user can be of two kinds: (i) Liar: a user who lies about his values, but otherwise follows the protocol correctly, or (ii) Breaker: a user who breaks the protocol, e.g., sends wrong messages. A malicious aggregator can break the protocol. In addition, it can collude with other malicious users.

To ensure privacy for users, we make a flexible assumption that at least a fraction of users (e.g., a majority) are honest. The lower bound h of the number of honest users is known a priori during deciding the noise generation parameters of the system. Remaining users and the aggregator can be arbitrarily malicious. (Similar assumption is generally made in cryptographic solutions.) The assumption is practical and flexible—for a user to know that his true values will remain private, he only needs to know that at least a certain fraction of other users are honest; he does not need to trust any single entity such as the aggregator.

On the other hand, to ensure a good utility guarantee for the aggregator, we assume that the aggregator queries a set of users that it generally trusts. Of the users the aggregator chooses to query, there can be at most l liars (l is small) and the remaining users are either honest or colluding/collaborating with the aggregator. There is fundamentally no way to ensure good utility if a large number of users lie about their data. The same is true if even a single user introduces an arbitrarily large noise. So we assume that there are no breakers: in practice, this can be arranged by ensuring that users' messages sent to the aggregator are generated and digitally signed by a trusted software implementing the protocol.⁶

In summary, our privacy guarantees hold even if a large number of users are malicious. This is crucial to make new privacy-aware users feel comfortable to join the system. Our utility guarantees hold if a small ($< l$) number of users lie and try to disrupt the final aggregate. This leaves the responsibility to the aggregator for choosing a good set of users to query. For example, if the aggregator can identify a malicious user (e.g., via some out-of-band mechanism), it can blacklist the user and exclude him from its queries. Our attack model is stronger than many previous techniques [12] that use the honest-but-curious model and disallow malicious agents.

Privacy Goals. We aim to enable the aggregator to estimate answers to aggregate queries. At the same time, an aggregator should not learn anything more, other than the aggregate answer, about honest individual users. Moreover, no user should learn anything about the values of other honest users, even if he colludes with the aggregator or other malicious users. We formalize our privacy requirement using differential privacy, which we discuss next.

2.2 Privacy Definition and Background

We now discuss differential privacy [5] that we use as our privacy definition and also review the standard perturbation technique [7] used for achieving differential privacy.

Differential Privacy. Informally, an algorithm is differentially private if its output is insensitive to changes in the data I_u of any single user u . This provides privacy because if similar databases, say differing in the data of one user, produce indistinguishable outputs, then the adversary cannot use the output to infer any single user's data. To formalize this notion, denote $I = I_1 \cup I_2 \dots \cup I_U$ the combined data from U users and $nbrs(I)$ the data obtained from adding/removing one user's data from I , i.e. $nbrs(I)$ consists of I' such that either $I' = I \cup I_u$ for $u \notin \{1, \dots, U\}$ or $I' = I - I_u$ for some $u \in \{1, \dots, U\}$.

DEFINITION 2.1 (ϵ -DIFFERENTIAL PRIVACY [5]). *Denote $A(I)$ the output of an algorithm A on input data I . Then A is ϵ -differentially private if all I, I' such that $I' \in nbrs(I)$, and any output x , the following holds:*

$$Pr[A(I) = x] \leq e^\epsilon Pr[A(I') = x]$$

where Pr is a probability distribution over the randomness of the algorithm.

Query Sensitivity. We now look at a standard technique proposed by Dwork et al. [7] for differential privacy. It can be used to answer any query whether it is just a single snapshot query, Q_i , or a recurring query sequence, $\mathbf{Q} = Q_1, \dots, Q_n$.

The technique works by adding random noise to the answers, where the noise distribution is carefully calibrated to the query.

⁵Our techniques can support arbitrary queries if run on a centralized server.

⁶Many software security systems rely on trusted software, for instance the clients' antivirus software, to work with untrusted clients [14].

The calibration depends on the query *sensitivity*—informally, the maximum amount the query answers can change given any change to a single user’s data I_u . If \mathbf{Q} is a query sequence, $\mathbf{Q}(I)$ and $\mathbf{Q}(I')$ are each vectors. Sensitivity then measures the distance between the two vectors. This is typically done using the L_1 distance metric, denoted as $|\mathbf{Q}(I) - \mathbf{Q}(I')|_1$, that measures the Manhattan distance $\sum_i |\mathbf{Q}_i(I) - \mathbf{Q}_i(I')|$ between these vectors. In this paper, we also use the L_2 distance metric, denoted as $|\mathbf{Q}(I) - \mathbf{Q}(I')|_2$ that measures the Euclidean distance $\sqrt{\sum_i (\mathbf{Q}_i(I) - \mathbf{Q}_i(I'))^2}$.

DEFINITION 2.2 (SENSITIVITY [7]). Let \mathbf{Q} be any query sequence. For $p \in \{1, 2\}$, the L_p sensitivity of \mathbf{Q} , denoted $\Delta_p(\mathbf{Q})$, is the smallest number such that for all I and $I' \in \text{nbrs}(I)$,

$$|\mathbf{Q}(I) - \mathbf{Q}(I')|_p \leq \Delta_p(\mathbf{Q})$$

For a single snapshot query Q_i , the L_1 and L_2 sensitivities are the same, and we write $\Delta(Q_i) = \Delta_1(Q_i) = \Delta_2(Q_i)$.

EXAMPLE 2.1. Consider a query Q counting the number of users whose weight in month 1 is greater than 200 lb. Then $\Delta(Q)$ is simply 1 as Q can differ by at most 1 on adding/removing a single user’s data. Now consider $\mathbf{Q} = Q_1, \dots, Q_n$, where Q_i counts users whose weight in month i is greater than 200 lb. Then $\Delta_1(\mathbf{Q})$ is n (for the pair I, I' which differ in a single user having weight > 200 in each month i) and $\Delta_2(\mathbf{Q}) = \sqrt{n}$ (for the same pair I, I').

Laplace Perturbation Algorithm (LPA). To guarantee differential privacy in presence of a trusted server, [7] proposes the LPA algorithm that adds suitably-chosen noise to the true answers. The noise is generated according to the Laplace distribution. Denote $Lap(\lambda)$ a random variable drawn from the Laplace distribution with PDF: $Pr(Lap(\lambda) = Z) = \frac{1}{2\lambda} e^{-\epsilon|Z|/\lambda}$. $Lap(\lambda)$ has mean 0 and variance $2\lambda^2$. Also denote $\mathbf{Lap}^n(\lambda)$ to be a vector of n independent $Lap(\lambda)$ random variables.

The LPA algorithm takes as input a query sequence \mathbf{Q} and parameter λ controlling the Laplace noise. LPA first computes the true answers, $\mathbf{Q}(I)$, exactly and then perturbs the answers by adding independent $Lap(\lambda)$ noise to each query answer in $\mathbf{Q}(I)$. More formally, it computes and outputs $\tilde{\mathbf{Q}} = \mathbf{Q}(I) + \mathbf{Lap}^n(\lambda)$. Differential privacy is guaranteed if the parameter λ of the Laplace noise is calibrated according to the L_1 sensitivity of \mathbf{Q} . The following theorem shown in [7] formalizes this intuition.

THEOREM 2.1 (PRIVACY [7]). $LPA(\mathbf{Q}, \lambda)$ is ϵ -differentially private for $\lambda = \Delta_1(\mathbf{Q})/\epsilon$.

Analyzing accuracy. To analyze the accuracy of the perturbed estimates returned by an algorithm, we quantify their error in the following way.

DEFINITION 2.3 (UTILITY). Let $A(\mathbf{Q})$ be an algorithm that returns a perturbed estimate $\tilde{\mathbf{Q}} = \tilde{Q}_1, \dots, \tilde{Q}_n$ for an input sequence $\mathbf{Q} = Q_1, \dots, Q_n$. We denote $error_i(A) = \mathbb{E}_A |\tilde{Q}_i - Q_i(I)|$ the expected error in the estimate of the i^{th} query Q_i . Here \mathbb{E}_A is the expectation taken over the possible randomness of A . To quantify error in the entire sequence $\tilde{\mathbf{Q}}$ returned by algorithm A , we use the standard notion of Root Mean Square Error (RMSE). We denote it as $error(A)$; it is computed as $\sqrt{\mathbb{E}_A \sum_{i=1}^n |\tilde{Q}_i - Q_i(I)|^2}$.

For example, $error_i(LPA) = \mathbb{E} |\tilde{Q}_i - Q_i| = \mathbb{E} |Lap(\lambda)| = \lambda$, while the RMSE, $error(LPA)$, is $\sqrt{\sum_{i=1}^n \lambda^2} = \sqrt{n}\lambda$.

Next we discuss the utility of the LPA algorithm while satisfying ϵ -differential privacy.

THEOREM 2.2 (UTILITY [7]). Suppose we fix $\lambda = \Delta_1(\mathbf{Q})/\epsilon$ so that $LPA(\mathbf{Q}, \lambda)$ is ϵ -differentially private. Then for all $i \in \{1, \dots, n\}$, $error_i(LPA) = \Delta_1(\mathbf{Q})/\epsilon$.

EXAMPLE 2.2. Recall the recurring query \mathbf{Q} of Eg. 2.1 that counts users having weight > 200 in each month $i = \{1, \dots, n\}$. Then $\Delta_1(\mathbf{Q}) = n$ and LPA gives an $error_i(LPA) = n/\epsilon$ in each query Q_i for ϵ -differential privacy. Also, the RMSE, $error(LPA)$, is $\sqrt{\sum_{i=1}^n n^2/\epsilon^2} = n^{3/2}/\epsilon$.

3. SKETCH OF THE PASTE SOLUTION

Before discussing PASTE, our solution for differential privacy over distributed time-series data, we provide an outline for it in this section. PASTE uses several existing primitives including Discrete Fourier Transform (DFT), homomorphic encryption (that allows aggregation of encrypted values without decryption), and threshold encryption (that requires a threshold number of users for decryption). We will review these techniques in Sections 4 and 5.

PASTE uses the following two protocol stages for answering a sequence \mathbf{Q} of n queries. The first stage is a method to improve the accuracy of query answers and is described using a trusted central server. The second stage is then used to obtain a distributed solution.

1) Fourier Perturbation Algorithm (FPA $_k$). To answer query sequence \mathbf{Q} of length n with small error under differential privacy, we design the FPA $_k$ algorithm. FPA $_k$ is based on compressing the answers, $\mathbf{Q}(I)$, of the query sequence using an orthonormal transformation. Intuitively, this means finding a k -length *compressed query sequence*, $\mathbf{F}^k = F_1^k, \dots, F_k^k$, where $k \ll n$, such that the answers, $\mathbf{F}^k(I)$, can be used to approximately compute $\mathbf{Q}(I)$. Then we can perturb $\mathbf{F}^k(I)$ instead of $\mathbf{Q}(I)$ using a lower noise (the noise actually reduces by a factor of n/k) while preserving differential privacy. An additional error creeps in since $\mathbf{F}^k(I)$ may not be able to reconstruct $\mathbf{Q}(I)$ exactly, but for the right choice of the compressed query sequence, \mathbf{F}^k , this reconstruction error is significantly lower than the perturbation error caused by adding noise directly to $\mathbf{Q}(I)$. A good \mathbf{F}^k can be found using any orthonormal transformation and we use the Discrete Fourier Transform (DFT) in our algorithm. The DFT compressed query sequence, \mathbf{F}^k , has large L_1 sensitivity and perturbing it for differential privacy requires new analysis that distinguishes our solution from other Fourier-based perturbation approaches [20]. We discuss this in detail in Sec. 4.

2) Distributed LPA (DLPA). To answer an aggregate-sum query sequence \mathbf{Q} distributedly under differential privacy, we propose the DLPA algorithm, a distributed version of the LPA algorithm discussed in Sec. 2.2. Our complete solution comprises of using FPA $_k$ for improving accuracy together with DLPA for distributed noise-addition. We describe how they are combined in a moment, but first we explain the DLPA algorithm. We explain DLPA for a single aggregate-sum query Q : the generalization to the sequence \mathbf{Q} is straight-forward and just requires n separate invocations, once for each query Q_i in \mathbf{Q} . Since Q is an aggregate-sum query, $Q(I) = \sum_{u=1}^U f_u(I_u)$ where the function f_u maps user u ’s data to numbers. Denote $x_u = f_u(I_u)$, so that $Q(I) = \sum_{u=1}^U x_u$.

The basic protocol is based on threshold homomorphic encryption and is shown in Fig. 2. To perturb $Q(I) = \sum_{u=1}^U x_u$, each user u adds a share of noise, n_u , to his private value x_u . To keep the estimation error small, the noise shares are chosen such that $\sum_{u=1}^U n_u$ is sufficient for differential privacy, but n_u alone is not sufficient: thus the value $x_u + n_u$ can not directly be sent to the aggregator. To address this, the user u computes encryptions of $x_u + n_u$ and sends it to the aggregator. Due to encryption, the

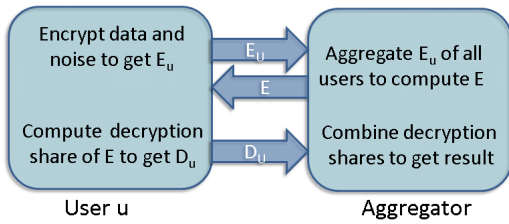


Figure 2: Basic Distributed Protocol (homomorphic property exploited to aggregate users’ encryption & threshold property to combine users’ decryption shares)

aggregator can not learn anything about x_u . However, using the *homomorphic* property, it can compute the encryption of the sum of all received values from all users thereby getting an encryption E of $\tilde{Q} = \sum_u (x_u + n_u)$. This encryption is then sent back to all users who use the *threshold* property to compute their respective decryption shares. Finally, the decryption shares are sent to the aggregator who combines them to compute the final decryption. The end result is a noisy differentially private estimate \tilde{Q} of $Q(I)$.

There are two challenges with the basic protocol described above. Firstly, the noise shares have to be generated in a way so that their sum is sufficient for differential privacy. Secondly, the aggregator can be malicious: the distributed protocol requires the aggregator to compute an encryption E of $\tilde{Q} = \sum_u (x_u + n_u)$ and send a decryption request to the users to help him decrypt E . But, the aggregator can cheat and request the decryption of wrong values, for instance, the encrypted private value of a single user, in which case the users will be inadvertently decrypting the private value of that user. We discuss how to solve these challenges in Sec. 5.

Putting the two together into PASTE. Now we explain how FPA_k and DLPA together give the complete solution in PASTE. To answer a query sequence \mathbf{Q} of n aggregate-sum queries, the aggregator first uses FPA_k to compute the k -length sequence \mathbf{F}^k . Due to the linearity of the DFT transformation, \mathbf{F}^k is another aggregate-sum query sequence. Now FPA_k requires to perturb the answers of the query sequence \mathbf{F}^k in order to get differential privacy. This is done by applying the DLPA algorithm on \mathbf{F}^k . The end result of DLPA is that the aggregator gets a noisy differentially estimate $\tilde{\mathbf{F}}^k$ of $\mathbf{F}^k(I)$. Then the aggregator computes the inverse DFT to reconstruct an estimate $\tilde{\mathbf{Q}}$ from $\tilde{\mathbf{F}}^k$. The final estimate $\tilde{\mathbf{Q}}$ has error characteristics of the FPA_k algorithm, but has been computed in a distributed way using DLPA.

We discuss FPA_k and DLPA in detail in next two sections.

4. FOURIER PERTURBATION ALGORITHM

We now describe in detail the FPA_k algorithm (the first stage of PASTE) for improving accuracy of query answers for long sequences. In this section, we assume a central trusted server: how to distribute the algorithm using DLPA (the second stage of PASTE) was briefly mentioned in the solution sketch and will be discussed in detail in Sec. 5. The FPA_k algorithm is based on the Discrete Fourier Transform, which we review briefly below.

4.1 The Discrete Fourier Transform

The DFT of a n -dimensional sequence \mathbf{X} is a linear transform giving another n -dimensional sequence, $\mathbf{DFT}(\mathbf{X})$, with j^{th} element given as: $\mathbf{DFT}(\mathbf{X})_j = \sum_{i=1}^n e^{2\pi\sqrt{-1}ji/n} \mathbf{X}_i$. Similarly one can compute the Inverse DFT as $\mathbf{IDFT}(\mathbf{X})_j = \frac{1}{n} \sum_{i=1}^n e^{-2\pi\sqrt{-1}ji/n} \mathbf{X}_i$. Furthermore, $\mathbf{IDFT}(\mathbf{DFT}(\mathbf{X})) = \mathbf{X}$.

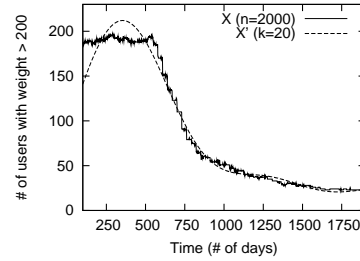


Figure 3: Reconstructed sequence \mathbf{X}' vs. original \mathbf{X}

Algorithm 4.1 FPA_k (Inputs: sequence \mathbf{Q} , parameter λ)

- 1: Compute $\mathbf{F}^k = \mathbf{DFT}^k(\mathbf{Q}(I))$.
- 2: Compute $\tilde{\mathbf{F}}^k = \text{LPA}(\mathbf{F}^k, \lambda)$
- 3: Return $\tilde{\mathbf{Q}} = \mathbf{IDFT}(\text{PAD}^n(\tilde{\mathbf{F}}^k))$

Denote $\mathbf{DFT}^k(\mathbf{X})$ as the first k elements of $\mathbf{DFT}(\mathbf{X})$. The elements of $\mathbf{DFT}^k(\mathbf{X})$ are called the Fourier coefficients of the k lowest frequencies and they compactly represent the high-level trends in \mathbf{X} . An approximation \mathbf{X}' to \mathbf{X} can be obtained from $\mathbf{DFT}^k(\mathbf{X})$ as follows: Denoting $\text{PAD}^n(\mathbf{DFT}^k(\mathbf{X}))$ the sequence of length n obtained by appending $n - k$ zeros to $\mathbf{DFT}^k(\mathbf{X})$, compute $\mathbf{X}' = \mathbf{IDFT}(\text{PAD}^n(\mathbf{DFT}^k(\mathbf{X})))$. Obviously \mathbf{X}' may be different from \mathbf{X} as ignoring the last $n - k$ Fourier coefficients may introduce some error. We denote $\text{RE}_i^k(\mathbf{X})$, short for reconstruction error at the i^{th} position, to be the value $|X'_i - X_i|$.

EXAMPLE 4.1. To give a sense of the reconstruction error, we consider a sequence \mathbf{X} of length $n = 2000$ representing the number of people with weight > 200 in a real dataset (more details in Section 7), counted once every day over 2000 days. Fig. 3 shows the reconstructed sequence, \mathbf{X}' , using $k = 20$ DFT coefficients along with the original sequence \mathbf{X} . \mathbf{X} shows the temporal trend in the # of overweight people in the dataset. As shown, \mathbf{X}' captures the trend accurately showing that the reconstruction error is small even when compressing from $n = 2000$ to $k = 20$ DFT coefficients.

4.2 The Algorithm

FPA_k is shown in Algorithm 4.1. Given a query sequence \mathbf{Q} , it begins by computing a compressed query sequence \mathbf{F}^k , comprising the first k Fourier coefficients in the DFT of $\mathbf{Q}(I)$. Then it perturbs \mathbf{F}^k using the LPA algorithm with parameter λ to compute a noisy estimate $\tilde{\mathbf{F}}^k$. This perturbation is done to guarantee differential privacy. Finally, the algorithm computes the inverse DFT of $\text{PAD}^n(\tilde{\mathbf{F}}^k)$ to get $\tilde{\mathbf{Q}}$, an approximation to the original query answers $\mathbf{Q}(I)$.

As with LPA, the parameter λ in FPA_k needs to be adjusted in order to get ϵ -differential privacy. Since FPA_k perturbs the sequence \mathbf{F}^k , λ has to be calibrated according to the L_1 sensitivity, $\Delta_1(\mathbf{F}^k)$, of \mathbf{F}^k . Analyzing this sensitivity is challenging as changing a single user’s data can affect the entire sequence, $\mathbf{Q}(I)$, resulting in large changes to each Fourier coefficient. Next we discuss how to bound this change, and obtain a value of λ that makes $FPA_k(\mathbf{Q}, \lambda)$ differentially private.

THEOREM 4.1. Denote $\mathbf{F}^k = \mathbf{DFT}^k(\mathbf{Q}(I))$ the first k DFT coefficients of $\mathbf{Q}(I)$. Then, (i) the L_1 sensitivity, $\Delta_1(\mathbf{F}^k)$, is at most \sqrt{k} times the L_2 sensitivity, $\Delta_2(\mathbf{Q})$, of \mathbf{Q} , and (ii) $FPA_k(\mathbf{Q}, \lambda)$ is ϵ -differentially private for $\lambda = \sqrt{k}\Delta_2(\mathbf{Q})/\epsilon$.

Proof (i) holds since $\Delta_2(\mathbf{F}^k) \leq \Delta_2(\mathbf{Q})$ (as the n Fourier coefficients have the same L_2 norm as \mathbf{Q} , while \mathbf{F}^k ignores the last $n - k$ Fourier coefficients), and $\Delta_1(\mathbf{F}^k) \leq \sqrt{k}\Delta_2(\mathbf{F}^k)$ (due to a standard inequality between the L_1 and L_2 norms of a sequence). (ii) follows since for $\lambda = \sqrt{k}\Delta_2(\mathbf{Q})/\epsilon \geq \Delta_1(\mathbf{F}^k)/\epsilon$, $\tilde{\mathbf{F}}^k = LPA(\mathbf{F}^k, \lambda)$ computed in Step 2 is ϵ -differentially private, and $\tilde{\mathbf{Q}}$ in step 3 is obtained using $\tilde{\mathbf{F}}^k$ only.

4.3 Analyzing accuracy

Example 2.2 gives an instance of a n -length query sequence for which LPA results in an error of n/ϵ to each query answer, making the answers useless for a large n . Intuitively speaking, FPA_k applies LPA on a length- k sequence. Hence the noise added during perturbation should be smaller. On the other hand, ignoring $n - k$ DFT coefficients in FPA_k results in an additional (but often much smaller) reconstruction error. Below we formalize this argument to compare the errors for the two algorithms in greater detail.

In general, LPA results in a large noise whenever the L_1 sensitivity of the query sequence is high. We define below irreducible query sequences that have the worst-possible L_1 sensitivity behavior.

DEFINITION 4.1 (IRREDUCIBLE QUERIES). *A query sequence $\mathbf{Q} = Q_1, \dots, Q_n$ is irreducible if its L_1 sensitivity, $\Delta_1(\mathbf{Q})$, is equal to the sum, $\sum_{i=1}^n \Delta(Q_i)$, of the sensitivities of its constituent queries Q_i 's.*

For all sequences \mathbf{Q} , $\Delta_1(\mathbf{Q}) \leq \sum_{i=1}^n \Delta(Q_i)$. Hence irreducible queries have the worst-possible L_1 sensitivities among all query sequences. Recurring queries over time-series data are often irreducible: one instance is the recurring query \mathbf{Q} of Example 2.2. The improvement of FPA_k over LPA is the most on irreducible query sequences (since LPA has the least accuracy for such queries). We compare the accuracies over irreducible queries below (note however that irreducibility is not required for FPA_k).

For simplicity, we assume W.L.O.G that the sensitivity $\Delta(Q_i) = 1$ for all $Q_i \in \mathbf{Q}$. If not, we can rescale the queries by defining the sequence $\mathbf{Q}' = \{Q'_1, \dots, Q'_n\}$ given as $Q'_i = Q_i/\Delta(Q_i)$. Then $\Delta(Q'_i) = \Delta(Q_i/\Delta(Q_i)) = \Delta(Q_i)/\Delta(Q_i) = 1$. Furthermore, $\mathbf{Q}(I)$ can be computed from $\mathbf{Q}'(I)$ by just multiplying with $\Delta(Q_i)$ at the i^{th} position. We call such sequences as normalized query sequences.

With this simplification, $\Delta_1(\mathbf{Q}) = n$ (as irreducibility means $\Delta_1(\mathbf{Q}) = \sum_{i=1}^n \Delta(Q_i) = \sum_{i=1}^n 1 = n$). Applying Theorem 2.2, we know that $error_i(\text{LPA}) = n/\epsilon$ in each query Q_i . On the other hand the following theorem shows the error of the FPA_k algorithm. Recall that $RE_i^k(\mathbf{Q}(I))$ is the reconstruction error for the i^{th} query, Q_i , caused by ignoring $n - k$ DFT coefficients of $\mathbf{Q}(I)$ and then computing the inverse DFT.

THEOREM 4.2. *Fix $\lambda = \sqrt{k}\Delta_2(\mathbf{Q})/\epsilon$ so that $FPA_k(\mathbf{Q}, \lambda)$ is ϵ -differentially private. Then for all $i \in \{1, \dots, n\}$, the $error_i(FPA_k)$ is $k/\epsilon + RE_i^k(\mathbf{Q}(I))$.*

Due to lack of space, the proof of the above theorem is deferred to the full version of the paper [22]. The theorem shows that the error by FPA_k for each query is $k/\epsilon + RE_i^k(\mathbf{Q}(I))$, while we know that LPA yields an error of n/ϵ . Since the reconstruction error, $RE_i^k(\mathbf{Q}(I))$, is often small even for $k \ll n$, we expect the error in FPA_k to be much smaller than in LPA. This hypothesis is confirmed in our experiments that show that FPA_k gives orders of magnitude improvement over LPA in terms of error.

Choosing the Right k . So far we have assumed that k is known to us. Since $error_i(FPA_k)$ is $k/\epsilon + RE_i^k(\mathbf{Q}(I))$, a good value of k

is important in obtaining a good trade-off between the perturbation error, k/ϵ , and the reconstruction error, $RE_i^k(\mathbf{Q}(I))$. If k is too big, the perturbation error becomes too big (giving the performance of LPA), while if k is too small the reconstruction error becomes too high.

We can often choose k based on prior assumptions about $\mathbf{Q}(I)$. For instance, if $\mathbf{Q}(I)$ is such that the Fourier coefficients corresponding to $\mathbf{Q}(I)$ decrease exponentially fast, then only a constant number (say $k=10$) of Fourier coefficients need to be retained during perturbation. Our experiments show that this naive method is applicable in many practical scenarios as Fourier coefficients of many real-word sequences decrease very rapidly [1].

However, for optimal performance, we need to adjust the value of k depending on the exact nature of $\mathbf{Q}(I)$. Computing k after looking at $\mathbf{Q}(I)$, however, compromises differential privacy. An algorithm to efficiently compute k in a differentially private way is described in Sec. 6.1.1.

5. DISTRIBUTED LPA (DLPA)

In both LPA and FPA_k , we assumed a trusted central server that stores the entire database I , computes the true answers $\mathbf{Q}(I)$, and adds noise for privacy. Next we discuss how to adapt these algorithms for the distributed setting, i.e. database $I = I_1 \cup I_2 \dots \cup I_U$ where I_u is the data of user u that she keeps to herself. We restrict ourselves to aggregate sum queries:⁷ i.e. \mathbf{Q} is such that $\mathbf{Q}(I) = \sum_{u=1}^U f_u(I_u)$ where f_u is an arbitrary function that maps user u 's database I_u to numbers.

We first construct the distributed LPA algorithm (DLPA for short) and then use it to construct distributed FPA_k . As discussed in Sec. 2, our algorithms guarantee privacy even against a malicious aggregator and malicious users as long as a majority of honest users exist. For the utility of the aggregate estimate to be good, we assume that among the malicious users, there are no breakers and at most l liars.

We explain DLPA over a single aggregate-sum query \mathbf{Q} : the generalization to query sequences is straight-forward and just requires multiple invocation of the algorithm, once for each query in the sequence. In this section, we also make a simplifying assumption that there are no communication failures: all links between the users and aggregator are maintained throughout the time required to answer the query. We relax this assumption and discuss fault tolerance in Sec. 6.2.

Let $\mathbf{Q} = \sum_{u=1}^U f_u(I_u)$ be the aggregate-sum query that we wish to answer. For simplicity of presentation, we assume f_u returns numbers in the set $\{1, 2, \dots, l\}$: our techniques work for any finite domain $\subseteq \mathbb{R}$. Denote $x_u = f_u(I_u)$. Then $\mathbf{Q} = \sum_{u=1}^U x_u$. For privacy, each user u adds a share of noise, n_u , to his private value x_u such that $\sum_{u=1}^U n_u$ is sufficient for differential privacy, but n_u itself is small so that total noise is low as possible. Hence, the value $x_u + n_u$ can not directly be sent to the aggregator. To address this, the user u computes encryptions of $x_u + n_u$ and sends it to the aggregator, which performs aggregation over encrypted data. We next review the cryptographic primitives that make this possible.

5.1 Basics: Encryption Scheme

DLPA is built upon the threshold Paillier cryptosystem [10]. The cryptosystem is set up by choosing an integer m such that (i) $m = pq$ where p and q are strong primes (i.e $p = 2p' + 1$ and $q = 2q' + 1$), and (ii) $\gcd(m, \phi(m)) = 1$. Once m is chosen, any number in \mathbb{Z}_m (the set $\{0, 1, \dots, m - 1\}$) can be encrypted. Also

⁷Note that FPA_k itself is more general and can support arbitrary queries when used in a centralized setting.

Algorithm 5.1 Encrypt-Sum(x_u, r_u)

- 1: User u generates a random $r_u \in \mathbb{Z}_m$ computes $c_u = \text{Enc}(x_u + r_u)$, and sends c_u to the aggregator.
 - 2: The aggregator computes $c = \prod_{u=1}^U c_u$.
-

denote \mathbb{Z}_m^* the subset of numbers in \mathbb{Z}_m that have a multiplicative inverse modulo m (eg. 0 does not have an inverse, but 1 has).

Key generation Choose a random element $\beta \in \mathbb{Z}_m^*$ and set $\lambda = \beta \times \text{lcm}(p, q)$. λ is the private key. Also set $g = (1+m)^a b^m \text{ mod } m^2$ for some randomly chosen $(a, b) \in \mathbb{Z}_m^* \times \mathbb{Z}_m^*$. The triplet (m, g, g^λ) forms the public key.

Encryption The encryption function Enc maps a plaintext message $t \in \mathbb{Z}_m$ to ciphertext $c \in \mathbb{Z}_{m^2}^*$. $\text{Enc}(t)$ is computed as $g^t r^m \text{ mod } m^2$ where $r \in \mathbb{Z}_m^*$ is a randomly chosen number.

Decryption Denote L the function $L(u) = (u - 1)/m$ for any $u = 1 \text{ mod } m$. The decryption of ciphertext $c \in \mathbb{Z}_{m^2}^*$ is the function $\text{Dec}(c) = \frac{L(c^\lambda \text{ mod } m^2)}{L(g^\lambda \text{ mod } m^2)}$

The encryption scheme has the following properties.

Homomorphic addition If c_i is a ciphertext for message t_i for $i \in 1, 2$, then $c_1 \cdot c_2$ is a ciphertext for message $t_1 + t_2$.

Distributed decryption Suppose the private key λ is shared by U users as $\lambda = \sum_u \lambda_u$ where λ_u is the private key for user u . Then decryption of a ciphertext c can be done distributedly (i.e. without any party knowing λ) as:

- Each user u computes his decryption share $c_u = c^{\lambda_u}$.
- The decryption shares are combined as $c' = \prod_{u=1}^U c_u$.
- Finally the decryption $t = \frac{L(c' \text{ mod } m^2)}{L(g^\lambda \text{ mod } m^2)}$ is computed.

5.2 Protocol for Computing Exact Sum

Let $Q = \sum_{u=1}^U x_u$ be the query to be computed by the aggregator. We first start with a protocol for computing Q at the aggregator exactly. We will subsequently enhance the protocol to include noise addition. As discussed in the solution sketch (Sec. 3), computing even exact sum is difficult against a malicious aggregator: instead of the aggregator computing the encrypted sum and sending it to the users for decryption, it can send false decryption requests (say the encrypted private value of a single user) and break privacy.

We use the threshold Paillier scheme for the protocol. In a key generation phase, the private key λ is generated and distributed among the users as $\lambda = \sum_{i=1}^U \lambda_i$. Thus the users all together can perform distributed decryption using their keys λ_u . Note that since the key generation needs to be done only once (irrespective of the # of queries to be answered), expensive secret-sharing protocols [24] can be used for this purpose.

The protocol executes in two phases. In the first phase, the aggregator computes the required Q in encrypted form. Then in the second phase, a distributed decryption protocol is run to recover Q from the encrypted form.

The first phase is shown in Algorithm 5.1. We call it Encrypt-Sum(x_u, r_u): each user u encrypts his private value, x_u , added to a randomly generated r_u . Note that r_u is known only to user u . The aggregator obtains all the encryptions and multiplies them to compute c . Due to the homomorphic properties of the encryption, the obtained c is an encryption of $\sum_{u=1}^U (x_u + r_u) = Q + \sum_{u=1}^U r_u$.

Since r_u 's are not known to the aggregator, decrypting c would not reveal any information about Q . However, the following modification of the distributed decryption protocol can be used to obtain Q exactly. Note that g^λ is publicly known.

Algorithm 5.2 Decrypt-Sum(c, r_u)

- 1: The aggregator sends c to each user u for decryption.
 - 2: User u computes decryption share $c'_u = c^{\lambda_u} g^{-r_u \lambda}$.
 - 3: The aggregator collects c'_u from each user, combines them to get $c' = \prod_{i=1}^U c'_i$, and computes the final decryption $\bar{Q} = \frac{L(c' \text{ mod } m^2)}{L(g^\lambda \text{ mod } m^2)}$
-

The above protocol is a minor modification of distributed decryption: user u multiplies an additional factor of $g^{-r_u \lambda}$ while generating his decryption share (step 2). We call this protocol Decrypt-sum(c, r_u). The following proposition shows the correctness of the protocol.

PROPOSITION 5.1. *Let $c = \text{Encrypt-sum}(x_u, r_u)$ and \bar{Q} be the decryption computed by $\text{Decrypt-sum}(c, r_u)$. Then $\bar{Q} = Q = \sum_{u=1}^U x_u$.*

Proof Sketch: The proof appears in the full version of the paper [22]. Here we give a sketch. As mentioned earlier, c obtained from $\text{Encrypt-sum}(x_u, r_u)$ is an encryption of $Q + \sum_{u=1}^U r_u$. Each user u corrects for his r_u in Step 2 of $\text{Decrypt-sum}(c, r_u)$. Thus the final \bar{Q} obtained is equal to Q .

Finally, we show that even though the Encrypt-sum and Decrypt-sum protocols can be used to compute the sum, $\sum_{u=1}^U x_u$, no other linear combinations can be computed. We saw that to compute the sum, the aggregator computed $c = \prod_u c_u$, where c_u is the encryption received by the aggregator from user u in the Encrypt-sum protocol (Step 2). Next we show that virtually no other encryption computed from these c_u 's can be decrypted in order to breach privacy.

THEOREM 5.1. *Suppose that the aggregator runs the $\text{Encrypt-sum}(x_u, r_u)$ protocol followed by the $\text{Decrypt-sum}(c', r_u)$ protocol for some c' of his choice. Let $c' = \prod_{u=1}^U c_u^{a_u}$ (where c_u are the encryptions sent by user u during Encrypt-sum protocol) such that $a_u - 1$ has an inverse mod m^2 (which implies $a_u \neq 1$) for some u . If the Decrypt-sum protocol decrypts c' correctly to give $\sum_{u=1}^U a_u x_u$, then there exists an attacker that breaks the security of the original distributed Paillier cryptosystem.*

Proof Sketch: The formal proof of security is quite involved and appears in the full version [22]: here we just highlight the intuition. Let $c = \text{Encrypt-sum}(x_u, r_u)$. Suppose the aggregator runs $\text{Decrypt-sum}(c', r_u)$ for $c' \neq c$, i.e. he sends wrong request for decryption, say the encryption of a single user's data (i.e., $a_u = 0$ for all users but one). In Step 2 of $\text{Decrypt-sum}(c', r_u)$ protocol, other users will wrongly be correcting for their random r_u 's that are not present in c' , making the decrypted value \bar{Q} completely random and useless.

5.3 Protocol for Computing Noisy Sum

Now we describe how to add noise in the distributed setting. As mentioned earlier, LPA requires us to compute $\tilde{Q} = Q + \text{Lap}(\lambda)$ where $\text{Lap}(\lambda)$ is a Laplace random variable with mean 0 and variance $2\lambda^2$. Denote $N(\mu, \sigma)$ to be a Gaussian random variable with mean μ and variance σ^2 . We shall generate Laplace noise using 4 Gaussian variables by exploiting the following property (proved in the full version [22]).

Algorithm 5.4 Encrypt-Noisy-Sum(x_u, r_u)

- 1: User u chooses five random numbers $r_u^1, r_u^2, \dots, r_u^5$ from \mathbb{Z}_m and computes $r_u = r_u^1 + r_u^2 - r_u^3 - r_u^4 + r_u^5$.
 - 2: User u generates four $N(0, \sqrt{2\lambda}/U)$ random variables y_u^1, \dots, y_u^4 .
 - 3: Let $c^j = \text{Encrypt-Sum-Squared}(y_u^j, r_u^j)$ for $j \in \{1, 2, 3, 4\}$.
 - 4: Let $c^5 = \text{Encrypt-Sum}(x_u, r_u^5)$.
 - 5: Aggregator computes $c = \frac{c^1 c^2 c^5}{c^3 c^4}$.
-

PROPOSITION 5.2. Let $Y_i \sim N(0, \lambda)$ for $i \in \{1, 2, 3, 4\}$ be four Gaussian random variables. Then $Z = Y_1^2 + Y_2^2 - Y_3^2 - Y_4^2$ is a $Lap(2\lambda^2)$ random variable.

The advantage of this decomposition is that Gaussian variables can be generated in a distributed fashion: To generate a $N(0, \lambda)$ variable, each user can generate a $N(0, \lambda/h)$ variable ($h = U/2$ is a lower bound on the number of honest users, i.e. a honest majority exists) and then the sum of these h , $N(0, \lambda/h)$, random variables gives the right $N(0, \lambda)$ variable. However, to compute a $Lap(\lambda)$ variable by Theorem 5.2, we need to compute squares of Gaussian random variables: for this we extend the Encrypt-Sum protocol described in the previous section to compute the encryption of $(\sum_{u=1}^U y_u)^2$ where y_u is the private $N(0, \lambda/h)$ of each user.

The protocol requires two randomly generated private keys $a_u, b_u \in \mathbb{Z}_m$ for each user u . The keys b_u are such that their sum for all users, $\sum_{u=1}^U b_u$, is 0. Denote a the sum $\sum_{u=1}^U a_u$. $Enc(a^2)$ is computed and made public in a key generation phase. The keys a_u, b_u need to be generated only once and expensive secret sharing protocols [24] can be used for this purpose. The protocol is shown below.

Algorithm 5.3 Encrypt-Sum-Squared(y_u, r_u) Protocol

- 1: User u computes $c_u = Enc(y_u + a_u + b_u)$ and sends it to the aggregator.
 - 2: The aggregator computes $c = \prod_{u=1}^U c_u$ and sends it to each user u .
 - 3: Each user u generates a random $r_u \in \mathbb{Z}_m$, computes $c_u = c^{y_u - a_u + b_u} Enc(r_u)$.
 - 4: The aggregator collects c_u from each user and computes $c' = (\prod_{u=1}^U c_u) Enc(a^2)$
-

We call the above protocol Encrypt-Sum-Squared(y_u, r_u). Due to the homomorphic properties of the encryption, the final c obtained in Encrypt-Sum-Squared(y_u, r_u) can be shown to be an encryption of $(\sum_u y_u)^2 + \sum_u r_u$.

Finally we can discuss the the noisy-sum protocol to add $Lap(\lambda)$ noise. The protocol is shown in Algorithm 5.4 and is called Encrypt-Noisy-Sum(x_u, r_u): each user generates 4 Gaussian, $N(0, \sqrt{2\lambda}/U)$, random variables in Step 2. Each of these 4 Gaussian variable is used to generate an encryption of a $N(0, \sqrt{\lambda}/2)$ Gaussian random variable in Step 3 using the Encrypt-Sum-Squared protocol. Then the user u generates an encryption for his private value x_u using the Encrypt-Sum protocol in Step 4. The random variables r_u^i are generated so as to force the right encryption to be computed by the aggregator in Step 5: the aggregator can only compute $c = \frac{c^1 c^2 c^5}{c^3 c^4}$; all others would have some r_u^i unbalanced. Due to the homomorphic properties of the encryption, the obtained c is an encryption of $\sum_u x_u + (\sum_u y_u^1)^2 + (\sum_u y_u^2)^2 - (\sum_u y_u^3)^2 - (\sum_u y_u^4)^2 + \sum_u r_u$. Also this c can be decrypted using decrypt-sum(c, r_u) protocol (Algorithm 5.2).

Next we state theorems showing the privacy and utility of the protocol. Due to space constraints, the proofs are deferred to the full version of the paper [22].

THEOREM 5.2 (PRIVACY). Let $c = \text{Encrypt-Noisy-Sum}(x_u, r_u)$ and $\tilde{Q} = \text{decrypt-sum}(c, r_u)$. If there are at least $U/2$ honest users, then $\tilde{Q} = Q + Lap(\lambda) + \text{Extra-Noise}$, where $Lap(\lambda)$ is the noise generated by honest users and the Extra-Noise is that generated by malicious users. Thus for $\lambda = \Delta(Q)/\epsilon$, ϵ -differential privacy is guaranteed independent of what the malicious users and aggregator choose to do.

THEOREM 5.3 (UTILITY). Let $c = \text{Encrypt-Noisy-Sum}(x_u, r_u)$ and $\tilde{Q} = \text{decrypt-sum}(c, r_u)$. If there are no malicious users, then $\tilde{Q} = Q + Lap(2\lambda)$. Finally, in presence of l malicious users that are all liars and no breakers, \tilde{Q} can deviate from $Q + Lap(2\lambda)$ by at most $l \times \Delta(Q)$.

PASTE. Now we explain how FPA_k and DLPA together give the complete solution in PASTE. Above we discussed how to compute the perturbed estimate \tilde{Q} for a single query Q . As mentioned earlier, extending distributed LPA for a n -length query sequence \mathbf{Q} is straightforward: apply the Encrypt-Noisy-Sum protocol n times, once for each $Q_i \in \mathbf{Q}$. This works since LPA consists of n independent perturbations for each of the n queries in the sequence.

Implementing FPA_k over the distributed setting is slightly more involved. Each user u first computes the answers of \mathbf{Q} over his data I_u . We denote the answers as $\mathbf{Q}(I_u) = Q_1(I_u), \dots, Q_n(I_u)$. Next the user u computes the first k DFT coefficients of the answers $\mathbf{Q}(I_u)$. Let us denote these k DFT coefficients by the sequence $\mathbf{F}^k(I_u)$. Recall that $\mathbf{F}^k(I)$ are the k DFT coefficients of $\mathbf{Q}(I)$, where $\mathbf{Q}(I)$ are the answers of \mathbf{Q} over the complete database I . By linearity of DFT, we know that $\mathbf{F}^k(I) = \sum_{u=1}^U \mathbf{F}^k(I_u)$: thus \mathbf{F}^k is another aggregate-sum query. Then distributed LPA can be used by the aggregator to compute the perturbed estimate $\tilde{\mathbf{F}}^k$ for $\mathbf{F}^k(I)$. Finally, the aggregator takes the inverse DFT transform of $\tilde{\mathbf{F}}^k$ to compute $\tilde{\mathbf{Q}}$, a perturbed estimate of \mathbf{Q} .

6. EXTENSIONS

We now describe two useful extensions of our algorithm. The first extension enables us to choose a good value of the parameter k (# of DFT coefficients used in FPA_k) if a central trusted server exists, while the second extension allows us to tolerate failures of users during the execution of the DLPA algorithm.

6.1 Choosing a Good Value of k

So far we have assumed that the value k is known to us before executing FPA_k . This was unavoidable in the distributed setting as each user needs to know k before computing his k DFT coefficients $\mathbf{F}^k(I_u)$. However, assuming a central server holds all the data I , then the server can compute k depending on the exact nature of $\mathbf{Q}(I)$. Computing k after looking at $\mathbf{Q}(I)$, however, compromises differential privacy. We present here an algorithm to efficiently compute k in a differentially private way. This algorithm can be used instead of LPA to accurately answer long query sequences in the centralized setting.

6.1.1 Differentially-private Sampling

Denote $\mathbf{F} = \text{DFT}(\mathbf{Q}(I))$, and \mathbf{F}^k (resp. \mathbf{F}^{n-k}) the first k (resp. last $n - k$) Fourier coefficients in \mathbf{F} . We want to choose a suitable k , and an estimate $\tilde{\mathbf{F}}^k$ for the first k Fourier coefficients in a differentially private way. One obvious method is to somehow compute k under differential privacy, and then use FPA_k to compute

Algorithm 6.1 SPA(Inputs: $\mathbf{Q}(I)$, parameter λ)

- 1: Compute $\mathbf{F} = \text{DFT}(\mathbf{Q}(I))$
 - 2: Denote $U'(k)$ the function $|\mathbf{F}^{n-k}|_2 + k\sqrt{n}/\epsilon$.
 - 3: Sample a $k \in \{1, \dots, n\}$ with probability $\propto e^{-U'(k)/\lambda}$.
 - 4: Let g be a $G(1/\lambda^2, (k+1)/2)$ random variable
 - 5: Compute $\tilde{\mathbf{F}}^k = \mathbf{F}^k + \mathbf{N}^k(0, \sqrt{g})$
 - 6: Compute $\tilde{\mathbf{Q}} = \text{IDFT}(\text{PAD}^n(\tilde{\mathbf{F}}^k))$.
-

$\tilde{\mathbf{F}}^k$. However this two step procedure results in additional noise, as computing k privately exhausts some privacy budget that could have been used for computing $\tilde{\mathbf{F}}^k$ using FPA_k . Next we discuss a differentially private sampling procedure that computes a suitable k along with $\tilde{\mathbf{F}}^k$ simultaneously.

THEOREM 6.1. Denote $U(k, \tilde{\mathbf{F}}^k)$ the function $|\mathbf{F}^{n-k}|_2 + |\tilde{\mathbf{F}}^k - \mathbf{F}^k|_2$. Then sampling k and $\tilde{\mathbf{F}}^k$ with probability proportional to $e^{-U(k, \tilde{\mathbf{F}}^k)/\lambda}$ satisfies ϵ -differential privacy for $\lambda = \sqrt{2}\Delta_2(\mathbf{Q})/\epsilon$.

Theorem 6.1 (proved in the full version [22]) shows a differentially private way of sampling k and $\tilde{\mathbf{F}}^k$. However, the sampling may not be efficient. This is because, even assuming that k has been sampled already, sampling $\tilde{\mathbf{F}}^k$ needs to be done with probability proportional to $e^{|\tilde{\mathbf{F}}^k - \mathbf{F}^k|_2/\lambda}$. In other words, $\tilde{\mathbf{F}}^k$ has to be sampled with probability based on its L_2 distance with another sequence \mathbf{F}^k . This is difficult as now elements of the sequence $\tilde{\mathbf{F}}^k$ cannot be sampled independently: if Pr is the sampling distribution, then $Pr(\tilde{F}_i^k = x | \tilde{F}_j^k = y) \neq Pr(\tilde{F}_i^k = x)$. On the other hand, in FPA_k , $\tilde{\mathbf{F}}^k$ was generated independently by adding $\text{Lap}^k(\lambda)$ random variables to \mathbf{F}^k . Nevertheless we discuss next an efficient way to sample from $U(k, \tilde{\mathbf{F}}^k)$.

6.1.2 Sampling Perturbation Algorithm (SPA)

Before giving the algorithm, we recall two kinds of random variables: (i) $N(\mu, \Sigma)$, that represents a normal random variable with mean μ and variance σ^2 (additionally, denote $\mathbf{N}^k(\mu, \sigma)$ a vector of k i.i.d normal variables), and (ii) $G(\theta, r)$, that represents a Gamma random variable with *PDF* given as $\frac{1}{\Gamma(r)\theta^r} x^{r-1} e^{-x/\theta}$, where $\theta > 0$, and $r > 0$ are parameters, and $\Gamma(\cdot)$ is the gamma function. This PDF is similar to the exponential distribution except for an extra factor of x^{r-1} .

Our sampling-based perturbation algorithm (SPA) is shown in Algorithm 6.1. In the first step it computes the entire DFT of \mathbf{Q} . Then in steps 2 and 3, it samples a value of k . Intuitively speaking, $U'(k) = |\mathbf{F}^{n-k}|_2 + k\sqrt{n}/\epsilon$ computed in the step 2 is the sum of the reconstruction error, $|\mathbf{F}^{n-k}|_2$ (this is the loss incurred by ignoring all elements in \mathbf{F}^{n-k}), and the perturbation error, $k\sqrt{n}/\epsilon$ (an additional factor of \sqrt{n} appears as this is the perturbation error for the whole sequence). In step 3, those values of k are more likely to be picked that give a lower $U'(k)$, i.e. give a better trade off between the reconstruction error and the perturbation error.

Once k has been sampled, the algorithm continues to sample $\tilde{\mathbf{F}}^k$. This is done by first picking a gamma random variable g in Step 4, and then perturbing \mathbf{F}^k by adding $\mathbf{N}^k(0, \sqrt{g})$ noise vector in Step 5. Even though, \mathbf{N}^k represents a vector of k independent normal variables, the $\tilde{\mathbf{F}}^k$ vector has not been generated in an independent fashion: g generated in Step 4 makes all elements of $\tilde{\mathbf{F}}^k$ correlated (for instance if g is picked to be 0, then all $\tilde{\mathbf{F}}^k = \mathbf{F}^k$). This makes sure that $\tilde{\mathbf{F}}^k$ has been generated in the right way, confirmed by the following theorem proved in the full version of the paper [22].

THEOREM 6.2 (PRIVACY). SPA(\mathbf{Q}, λ) is ϵ -differentially private for $\lambda = \sqrt{2}\Delta_2(\mathbf{Q})/\epsilon$.

Finally, we show that it always makes sense to run our SPA algorithm as opposed to LPA. If \mathbf{Q} is compressible, SPA will sample a good $k \ll n$ decreasing the error significantly (as confirmed in our experiments). However, if \mathbf{Q} is not compressible, the following theorem (proved in the full version [22]) shows that no matter what the total error of SPA would be at most a factor $\sqrt{2} \log n$ times worse than LPA for any normalized irreducible query sequence⁸ \mathbf{Q} .

THEOREM 6.3 (UTILITY). Let \mathbf{Q} be any normalized irreducible query sequence. Fix $\lambda_1 = \Delta_1(\mathbf{Q})/\epsilon$ and $\lambda_2 = \sqrt{2}\Delta_2(\mathbf{Q})/\epsilon$ such that LPA(\mathbf{Q}, λ_1) and $\text{FPA}_k(\mathbf{Q}, \lambda_2)$ are ϵ -differentially private. Then $\text{error}(\text{SPA}) \leq \sqrt{2} \cdot (\log n) \cdot \text{error}(\text{LPA})$.

6.2 Fault-tolerance

Our distributed protocol has two phases. In the Encrypt-Noise-Sum(x_u, r_u) phase, the aggregator computes an encryption c for the noisy sum of the private values, x_u , and the random values, r_u , of all users. Then in the Decrypt-Sum phase, all users need to correct for their respective random values r_u . In addition, since the secret key λ is shared as the sum $\sum_{u=1}^U \lambda_u$ of the private keys λ_u 's, all users need to send their decryption shares in order to decrypt. This makes the protocol susceptible to failures: if a single user does not respond in the decrypt-sum phase, no decryption can be obtained.

The solution is to (i) use a (T, U) -threshold decryption [10] scheme in which any T valid decryption shares out of U possible ones can be used to decrypt, and (ii) instead of choosing a completely random $r_u \in \mathbb{Z}_m$ during the encrypt-sum protocol, a user chooses $r_u \sim \text{Lap}(\Delta(\mathbf{Q})/\epsilon)$, i.e. r_u is chosen from the Laplace distribution sufficient for guaranteeing differential privacy. This r_u is sufficient: to minimize noise the aggregator has an incentive for adding all users' data when computing the aggregate encryption to be sent back for decryption, and leaving out a single user's data results in a Laplace noise sufficient for differential privacy.

Having seen the extensions still ensure privacy, let us see how they help in case of failures. Firstly, if f users fail to send their decryption shares, the (T, U) -threshold property ensures that a decryption can still be computed as long as $f < U - T$. Furthermore, not using the decryption share of f users means that the random value r_u of each of the f users is left uncorrected for, at the end of decryption. This results in an extra noise of $f \text{Lap}(\Delta(\mathbf{Q})/\epsilon)$ variables. It can be shown that the expected sum of these f variables increases as $\sqrt{f} \Delta(\mathbf{Q})/\epsilon$. In other words, the noise increases as the square root of the number of user failures f as demonstrated in our experiments. As long as f is small, this behavior is acceptable. For a large f , the distributed protocol for computation of \mathbf{Q} has to be repeated.

Note that in the above solution, we are concerned about failures of users that happen after the first phase and before the second phase of our protocol. Failures that happen before or after the complete execution of the protocol for a query do not affect the accuracy of our protocol. Since execution of our protocol for a single query takes less than a few seconds in practice, failures within this small time window is rare and are of small size. In Section 7.4, we experimentally show that the impact of such failures is small.

⁸Recall that normalized query sequences have individual query sensitivities rescaled to 1 while irreducible sequences have L_1 sensitivity equal to the sum of individual query sensitivities.

7. EXPERIMENTS

We have implemented PASTE in Java, using the BigInteger library for cryptographic components of the algorithms. This section evaluates our prototype using a 2.8 GHz Intel Pentium PC with 1GB RAM.

Data sets. We use three real data sets in our evaluation.

- **GPS:** This is the GPS trace from Microsoft’s Multiperson Local Survey (MLS) project [17]. The GPS trace was collected from 253 voluntary drivers, mostly in the Seattle, Washington area, covering about 135,000 kilometers of driving. The entire trace has approximately 2.3 million timestamped latitude/longitude points comprising about 16,000 discrete trips. The median interval between recorded points on a trip is 6 seconds and 62 meters. Thus, successive points on a trip are highly correlated.
- **Weight:** This is trace of body weights, collected from an online weight-monitoring website⁹. The trace contains daily weight data of about 300 users for a period of up to 5 years.
- **Traffic:** This data, collected from Department of Transportation of San Antonio, Texas¹⁰, reports volume and speed data at about 30 intersections in the city. We use a 6-months long trace, where data is reported once every 15 minutes.

Queries. For evaluating utility, we consider the following 6 query sequences. The first two queries are on the **GPS** data set, the next two queries are on the **Weight** data set, and the last two queries are on the **Traffic** dataset.

- Query *G1*: A query sequence whose value at each timestamp is an histogram over locations counting the number of people at each location. The locations are obtained by dividing the map into a 50×50 grid. The query sequence has 2000 different timestamps spread uniformly over 2 weeks.
- Query *G2*: The number of people at a fixed location evaluated at 2000 different timestamps spread uniformly over 2 weeks.
- Query *W1*: Number of people with weight > 200 lb on each day for a total of 2000 days.
- Query *W2*: A query sequence whose value at each timestamp is the correlation coefficient between two month-long time-series: (i) people with weight > 200 on day i of the month, and (ii) people with weight decreasing on day i of the month.
- Query *T1* and Query *T2*: Occupancy at two different city intersections with an interval of 15 minutes for a total of 6 months.

As the examples suggest, our algorithm can support a wide variety of queries. Unless otherwise stated, we use $k = 30$ for the FPA_k algorithm and 1-differential privacy ($\epsilon = 1$) as the privacy requirement in our experiments.

7.1 Accuracy of Answers

We first evaluate the accuracy of outputs of our FPA_k algorithm and compare it with the LPA algorithm. We report error percentage of a sequence of output estimate, which is the total L_2 error in the estimate (See Def. 2.3) normalized by the maximum possible L_2 value of a query answer.

Figure 4 reports the average error percentage of FPA_k and LPA algorithms for the 6 query sequences. Average error percentage is computed as the average over 100 runs of each algorithm, the variance in the error percentages are represented by the error bars.

⁹<http://www.hackersdata.com>

¹⁰<http://www.transguide.dot.state.tx.us/>

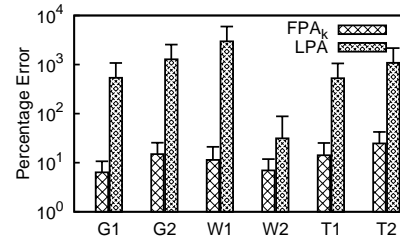


Figure 4: Average errors of FPA_k and LPA

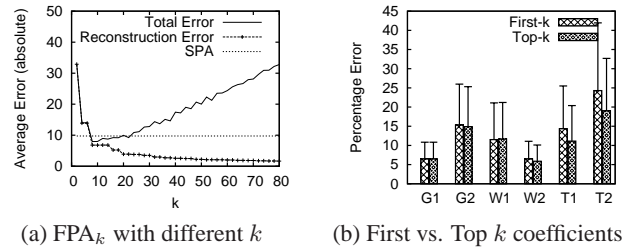


Figure 6: Effect of DFT parameters

As shown, FPA_k has orders of magnitude better error percentage than LPA (graph shows y-axis using a log-scale). In fact, LPA has an error percentage $\gg 100\%$ showing that estimates obtained are completely meaningless.¹¹ On the other hand, FPA_k has error percentages of $< 20\%$ that are mostly acceptable in practice.

Figure 5 demonstrates the errors of both algorithms more visually. Fig 5(a) and (c) plot outputs of the FPA_k algorithm along with the exact answers of *W1* and *G2*. Fig 5(b) and (d) plot the estimates obtained by LPA for the same two queries for the same privacy parameter. The graphs show that FPA_k estimates follow the curve of the exact answer very closely, while LPA results in estimates that are practically meaningless.

7.2 Effect of DFT Parameters

Value of k . To understand the impact of k on the accuracy of FPA_k , we vary the value of k in evaluating the query *W1*. Fig. 6(a) demonstrates the results. It shows that as k increases the total error decreases at first, then reaches a minimum, and finally increases. This is because the total error of FPA_k is a combination of reconstruction error and perturbation error. Reconstruction error decreases with increasing values of k (also shown in Figure 6(a)), while perturbation error increases with increasing k . The total error is minimized for an optimal value of k that finds a sweet spot between these two types of errors. This highlights the importance of choosing a good value of k for FPA_k .

In the distributed case, there is no way to choose the best k by looking at query answers. Thus k has to be predetermined, which is the reason we used a fixed value of $k = 30$ in our previous experiments (and it worked reasonably well for all queries). However the graph shows that for the *W1* query, $k = 10$ would have given a better trade-off, and the results of FPA_k would have been better if we used $k = 10$. In a centralized setting, our sampling perturbation algorithm (SPA) can be used to privately sample the value of

¹¹Laplace noise can send query answers outside their range making the error $> 100\%$. Truncation at range boundaries can make error = 100%, but is not done here to reflect the true magnitude of the noise.

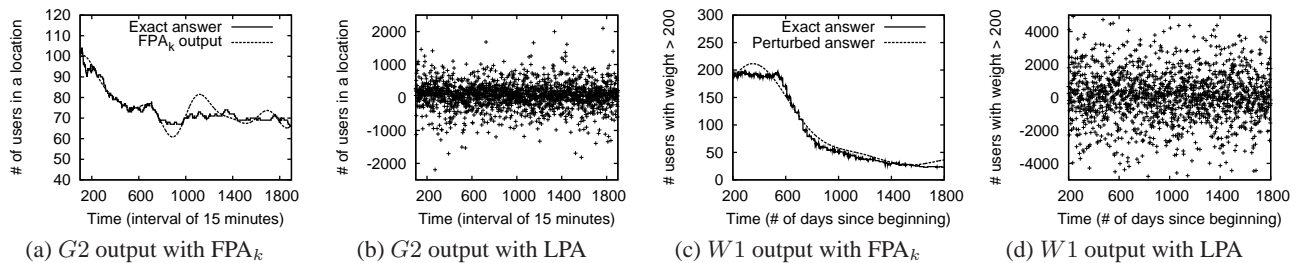


Figure 5: Point-wise errors of FPA_k and LPA over time

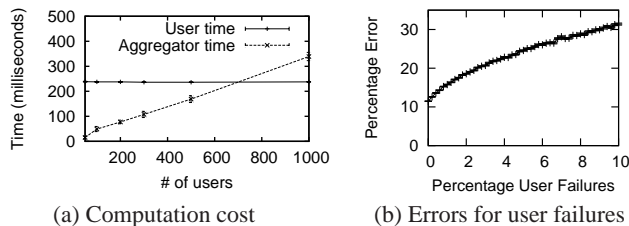


Figure 7: Evaluation of Distributed LPA

k . The error for that algorithm is also shown in the figure by the horizontal line and it is quite near the optimal error for $k = 10$.

First- k vs. Top- k DFT coefficients. Our algorithms, both FPA_k and SPA, choose the first k DFT coefficients, for a given value of k . The choice of leading coefficients is unavoidable in a distributed setting since all the users need to pick the same k DFT coefficients for our algorithms to work. Even in the presence of the centralized server, choosing the best set of k DFT coefficients in a differentially private way is inefficient. However, if possible, it is always better to choose the k largest DFT coefficients (i.e., top- k coefficients), since they give lower reconstruction error than the first- k coefficients. This leads to the natural question: how much do we sacrifice in accuracy for using the first- k DFT coefficients (which gives us differential privacy) instead of using the top- k coefficients (which does not give differential privacy)?

Figure 6(b) answers the above question. It shows the errors of FPA_k for different queries and compares it with the errors of a hypothetical algorithm that uses the top- k DFT coefficients (but may not necessarily be differentially private). Both the algorithms use $k = 30$ coefficients. The graph shows (with a linear scale for y -axis) that we do not lose much even if we just pick the first k coefficients (indicating that they generally are the largest coefficients). A substantial difference occurs only in $T1$ and $T2$ queries since they have periodic behavior at slightly higher frequencies indicating there largest coefficients are not the first k . Even for these queries, the difference is less than 5%.

7.3 Computation & Communication Overhead

Fig. 7(a) shows the computational overhead of our algorithm for computing a single aggregate-sum query. The graph shows the computation times (averaged over 100 runs) at a user and at the aggregator, as a function of the number of users involved in the query. The average computation time at a user is independent of the number of users and remains nearly constant. On the other hand, the time required by the aggregator increases linearly with the number of users. In all cases, the computation overhead is quite small, most of which is spent for cryptographic operations. We have also mea-

sured the communication overhead of our prototype and found that the algorithm has a small overhead: 0.5 Kb for each user (considering both incoming and outgoing messages) and 0.5 Kb times the # of users for the aggregator.

7.4 Effect of Failures

Fig. 7(b) shows the fault-tolerance behavior of our algorithm based on the extensions described in Sec. 6.2. The fault-tolerant algorithm is implemented and used to answer the $W1$ query of length $n = 2000$ using $k = 30$ DFT coefficients. We again report the percentage error: the total L_2 error in the computed estimate normalized by the maximum L_2 value of \mathbf{Q} . The percentage user failure is the percentage of users who fail during the protocol.

The graph shows that the error increases with the square root of the number of users failing and is quite reasonable with even say 5% failure rate. Note that the accuracy of our algorithm is affected only by the failures that happen *during* execution of the algorithm to answer a query, which is typically less than a few seconds in a distributed system. Failure within this short time should be rare and small in size in practice. In a very rare occasion, if a large number of users fail during this small time window, the algorithm should be started from the beginning.

8. RELATED WORK

Relational Data. Many approaches for relational data have been proposed that support formal definitions of privacy such as differential privacy [7, 15, 27], $\rho_1\rho_2$ breach [9], adversarial privacy [21, 23], etc. Among these, most relevant to our work are those that focus on query answering [7, 15, 27]. [7] proposes the LPA algorithm that adds independent Laplace noise to each query answer. The noise is calibrated according to the L_1 sensitivity. Recurring query sequences over time-series data have high L_1 sensitivity ($O(n)$ for a n length sequence) and thus LPA does not accurately answer such queries.

To improve accuracy of query answers under differential privacy, [27] focuses on range-count queries, i.e. count queries where the predicate on each attribute is a range. The main idea is to take a sequence of range-count queries that are disjoint and then perturb the entire Discrete Wavelet Transform of the sequence. Such disjoint range-count queries have a L_1 sensitivity of $O(1)$ and can be accurately answered using their technique. However the technique can not be used to answer recurring sequences over time-series data accurately, owing to their high L_1 sensitivity.

Another approach for answering query sequences accurately is proposed in [15], where constraints over the answers of multiple queries are exploited to improve accuracy. Usually no constraints exist over recurring query sequences and thus the method is largely inapplicable for time-series data. Similarly, [2] uses Fourier analysis to enforce consistency in query answers. Their Fourier analysis

uses a different perturbation algorithm whose goal is to compute consistent solutions efficiently, which is different from our goal of compressing time-series data.

Distributed protocols. Many Secure Multiparty Computation (SMC) techniques have been proposed to evaluate functions securely over data from multiple users [12]. In general, such techniques are slow [28] and infeasible for practical scenarios. However, if the function to be evaluated is a sum function, then efficient techniques using threshold homomorphic functions have been proposed [4, 13]. Such techniques often require a communication mechanism with a broadcast facility, or a mechanism enabling a user to verify the computation done by other users. This is not possible in our setting where the aggregator does all the computations and the users do not have enough resources to check whether those computations are correct. In addition, we need to compute noisy sum (in order to guarantee differential privacy) which is more difficult than computing just the sum of the inputs. To the best of our knowledge, [6] is the only known technique that computes noisy sum. However, it uses expensive secret-sharing protocols leading to a computation load of $O(U)$ per user, where U is the number of users. This makes the technique infeasible for large U .

Techniques for Time-series data. Most work on time-series data assume the centralized setting: a trusted server publishes an anonymized version of the data of all users and aggregate queries can then be run on the published data. [11] publishes data by adding a virtual user whose data is sampled from a public noise distribution. [20] works on a single time-series data by obtaining an orthonormal transform (such as DFT), then adding noise only to large coefficients of the transform, and finally obtaining the inverse transform. The main difference between our technique and theirs is the lack of formal privacy guarantee. Techniques [11, 20] show privacy by protecting against specific attacks (such as linear least-square regression or linear filtering in [20]), however no formal privacy guarantee (such as differential privacy) is provided. Furthermore no distributed solution is discussed. Same holds for most works on location data privacy (See [18] and the references therein).

9. CONCLUSION

We have proposed PASTE, a suite of novel algorithms to privately answer queries on distributed time-series data. PASTE uses the FPA_k algorithm to answer long query sequences over correlated time series data in a differentially private way. FPA_k perturbs k DFT coefficients of an answer sequence, thereby improving the accuracy for an n -length query sequence from $\Theta(n)$ of existing algorithms to roughly $\Theta(k)$, if the k DFT coefficients can accurately reconstruct all the query answers. For achieving differential privacy in distributed setting, PASTE uses the DLPA algorithm that implements Laplace noise addition in a distributed way with $O(1)$ complexity per user. Our experiments with three real data sets show that PASTE improves accuracy of query answers by orders of magnitude and also scales well with a large number of users.

Acknowledgments. We would like to thank Josh Benaloh, Melissa Chase, Seny Kamara, Frank McSherry, and Mariana Raykova for many helpful discussions. Eric Horvitz, John Krumm, and Paul Newson kindly gave us access to the GPS dataset.

10. REFERENCES

- [1] AGRAWAL, R., FALOUTSOS, C., AND SWAMI, A. N. Efficient similarity search in sequence databases. In *FODO* (1993).
- [2] BARAK, B., CHAUDHURI, K., DWORK, C., KALE, S., MCSHERRY, F., AND TALWAR, K. Privacy, accuracy, and consistency too: a holistic solution to contingency table release. In *PODS* (2007).
- [3] BLUM, A., DWORK, C., MCSHERRY, F., AND NISSIM, K. Practical privacy: The suLQ framework. In *PODS* (2005).
- [4] CRAMER, R., DAMGARD, I., AND NIELSEN, J. B. Multiparty computation from threshold homomorphic encryption. In *EUROCRYPT* (2001).
- [5] DWORK, C. Differential privacy: A survey of results. In *TAMC* (2008).
- [6] DWORK, C., KENTHAPADI, K., MCSHERRY, F., MIRONOV, I., AND NAOR, M. Our data, ourselves: Privacy via distributed noise generation. In *EUROCRYPT* (2006).
- [7] DWORK, C., MCSHERRY, F., NISSIM, K., AND SMITH, A. Calibrating noise to sensitivity in private data analysis. In *TCC* (2006).
- [8] EISENMAN, S. B., MILUZZO, E., LANE, N. D., PETERSON, R. A., AHN, G.-S., AND CAMPBELL, A. T. The bikenet mobile sensing system for cyclist experience mapping. In *ACM SenSys* (2007).
- [9] EVFIMIEVSKI, A., GEHRKE, J., AND SRIKANT, R. Limiting privacy breaches in privacy preserving data mining. In *PODS* (2003).
- [10] FOUQUE, P., POUPARD, G., AND STERN, J. Sharing decryption in the context of voting or lotteries. In *FC: International Conference on Financial Cryptography* (2000), LNCS, Springer-Verlag.
- [11] GANTI, R. K., PHAM, N., TSAI, Y.-E., AND ABDELZAHER, T. F. PoolView: stream privacy for grassroots participatory sensing. In *ACM SenSys* (2008).
- [12] GOLDREICH, O. Secure multi-party computation.
- [13] GRITZALIS, D. Secure electronic voting.
- [14] GUHA, S., REZNICHENKO, A., HADDADI, H., AND FRANCIS, P. Serving ads from localhost for performance, privacy, and profit. In *HotNets* (2009).
- [15] HAY, M., RASTOGI, V., MIKLAU, G., AND SUCIU, D. Boosting the accuracy of differentially-private histograms through consistency. In *VLDB 2010*.
- [16] HULL, B., BYCHKOVSKY, V., ZHANG, Y., CHEN, K., GORACZKO, M., MIU, A., SHIH, E., BALAKRISHNAN, H., AND MADDEN, S. Cartel: a distributed mobile sensor computing system. In *ACM SenSys* (2006).
- [17] KRUMM, J. A survey of computational location privacy. *Personal and Ubiquitous Computing*, 2008 (2008).
- [18] KRUMM, J., AND HORVITZ, E. Predestination: Where do you want to go today? *IEEE Computer* (2007).
- [19] MACHANAVAJHALA, A., GEHRKE, J., KIFER, D., AND VENKITASUBRAMANIAM, M. l-diversity: Privacy beyond k-anonymity. In *ICDE* (2006).
- [20] PAPADIMITRIOU, S., LI, F., KOLLIOS, G., AND YU, P. S. Time series compressibility and privacy. In *VLDB* (2007).
- [21] RASTOGI, V., HAY, M., MIKLAU, G., AND SUCIU, D. Relationship privacy: Output perturbation for queries with joins. In *PODS 2009*.
- [22] RASTOGI, V., AND NATH, S. Differentially private aggregation of distributed time-series with transformation and encryption. Tech. Rep. MSR-TR-2009-186, Microsoft Research, 2009.
- [23] RASTOGI, V., SUCIU, D., AND HONG, S. The boundary between privacy and utility in data publishing. In *VLDB* (2007).
- [24] SHAMIR, A. How to share a secret. *Communications of the ACM* 22, 11 (1979).
- [25] SHILTON, KATIE, B. J., ESTRIN, D., GOVINDAN, R., , AND KANG, J. Designing the personal data stream: Enabling participatory privacy in mobile personal sensing. In *37th Research Conference on Communication, Information and Internet Policy (TPRC)* (2009).
- [26] SWEENE, L. k-anonymity: A model for protecting privacy. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems* (2002).
- [27] XIAO, X., WANG, G., AND GEHRKE, J. Differential privacy via wavelet transforms. In *ICDE 2010*.
- [28] YAO, A. C.-C. Protocols for secure computations (extended abstract). In *FOCS* (1982).