# Continuous Release of Data Streams under both Centralized and Local Differential Privacy

Tianhao Wang<sup>1</sup>, Joann Qiongna Chen<sup>2</sup>,
Zhikun Zhang<sup>3</sup>, Dong Su<sup>4</sup>, Yueqiang Cheng<sup>5</sup>, Zhou Li<sup>2</sup>, Ninghui Li<sup>1</sup>, Somesh Jha<sup>6</sup>

<sup>1</sup>Purdue University, <sup>2</sup>University of California, Irvine,
<sup>3</sup>CYOPA Mala Lie Garage State Character of Mala Lie Garage State Charac

<sup>3</sup>CISPA Helmholtz Center for Information Security, <sup>4</sup>Alibaba, <sup>5</sup>Baidu Security, <sup>6</sup>University of Wisconsin, Madison {tianhaowang, ninghui}@purdue.edu, {qiongnac, zhou.li}@@uci.edu,

zhikun.zhang@cispa.saarland, sudong.tom@gmail.com, chengyueqiang@baidu.com, jha@cs.wisc.edu

Abstract—In this paper, we study the problem of publishing a stream of real-valued data satisfying differential privacy (DP). One major challenge is that the maximal possible value can be quite large; thus it is necessary to estimate a threshold so that numbers above it are truncated to reduce the amount of noise that is required to all the data. The estimation must be done based on the data in a private fashion. We develop such a method that uses the Exponential Mechanism with a quality function that approximates well the utility goal while maintaining a low sensitivity. Given the threshold, we then propose a novel online hierarchical method and several post-processing techniques.

Building on these ideas, we formalize the steps into a framework for private publishing of stream data. Our framework consists of three components: a threshold optimizer that privately estimates the threshold, a perturber that adds calibrated noises to the stream, and a smoother that improves the result using post-processing. Within our framework, we design an algorithm satisfying the more stringent setting of DP called local DP (LDP). To our knowledge, this is the first LDP algorithm for publishing streaming data. Using four real-world datasets, we demonstrate that our mechanism outperforms the state-of-the-art by a factor of 6-10 orders of magnitude in terms of utility (measured by the mean squared error of answering a random range query).

## I. Introduction

Continuous observation over data streams has been utilized in several real-world applications. For example, security companies continuously analyze network traffic to detect abnormal Internet behaviors [10]. Despite tangible benefits, analyzing and releasing data raises privacy concerns when these data contain sensitive individual information. Directly publishing raw statistics may reveal individual users' private information, such as daily routines and living habits. For instance, electricity usage data from smart meters can reveal whether a user is at home or even what household appliances are in use at some specific time [31]. These potential privacy risks would discourage data sharing.

A promising technique for private statistics releasing is differential privacy (DP) [18], which has become the gold standard in the privacy-research community. Informally, any

when any single record in the dataset is replaced. The closeness is quantified by a parameter  $\epsilon$ . A smaller  $\epsilon$  means a better privacy guarantee.

To provide statistics (e.g., the moving average) over data

algorithm satisfying DP has the property that its output distri-

bution on a given database is close to the output distribution

streams while satisfying DP, a widely-accepted technique is hierarchical algorithm [9], [19]. The idea is to partition the time series into multiple granularities and add noise proportional to the upper bound of the data. In [33], the authors observed that in many cases, data in the stream is concentrated below a threshold, which is much smaller than the upper bound. A commonly used technique in DP is called contribution limitation, which truncates the values using a specified threshold  $\theta$  (i.e., values larger than  $\theta$  will be replaced by  $\theta$ ) e.g., [47], [28], [15], [45], [27]). The rationale is to reduce the noise (now the DP noise is proportional only to the threshold) while maintaining most of the utility in the data. To find such a threshold while maintaining DP, [33] developed a method based on smooth sensitivity [32]. The result can then be applied to the hierarchical algorithm to publish streams with improved utility.

We find three key limitations in [33]. First, it tries to privately find the 99.5-th percentile to serve as the threshold  $\theta$ . Unfortunately, using the 99.5-th percentile (or any other fixed percentile) is unlikely to work across different settings with varying  $\epsilon$  values, data ranges, and data distributions. Second, in order to achieve an analytical upperbound on the error caused by truncation (this error is also called bias), it further increases the estimated 99.5-th percentile by first adding a positive bias term, and then multiplying a coefficient greater than 1. As a result, the chosen  $\theta$  value is often unnecessarily large. When  $\epsilon$  is small (e.g.,  $\epsilon \leq 0.1$ ), the chosen  $\theta$  values are usually larger than the maximal possible value, defeating the purpose of choosing the threshold in the first place. Third, the method utilizes a basic hierarchical approach directly to output the stream, and does not fully take advantage of possible post-processing optimizations. As a result, the accuracy of the outputted stream is quite low, and the results are even worse when answering range queries with small selectivity.

In this paper, we propose a new approach by addressing the above limitations. First, we design a data-dependent method to find the threshold  $\theta$  by looking into the data distribution (with DP guarantee), without the need to estimate a fixed percentile. Second, we identify that the reason for the conflict between

the threshold  $\theta$  and the small  $\epsilon$  is that the existing method only considers the bias error due to truncation, and ignores the noise error of DP. To address this conflict, we propose an algorithm built upon the differentially private exponential mechanism (EM [30]) by considering both errors simultaneously. The new algorithm is able to select an appropriate  $\theta$  with minimal overall errors. Third, we propose a new hierarchical algorithm to get accurate results. Specifically, we design an on-line consistency algorithm to enforce the consistency over the noisy estimates on the hierarchy to provide better utility. In addition, we observe that the estimates in the lower levels of the hierarchy tend to be overwhelmed by the noise, leading to a low signal-noise ratio. Thus, we further extend the algorithm to prune the lower-level nodes based on optimization criteria. Our new hierarchical algorithm is also able to handle infinite streams

We generalize the procedure of the above algorithms into a new framework for streaming data publication. It consists of three components: a Threshold optimizer, a Perturber, and a Smoother. The threshold optimizer consumes a portion of the input stream, and finds a threshold  $\theta$ . It then truncates any incoming value by  $\theta$  and sends them to the perturber. The perturber adds noise to each incoming element of the stream, and releases noisy counts to the smoother. Finally, the smoother performs further post-processing on the noisy counts and outputs the final stream. Integrate with the new algorithms described above, we call our solution ToPS.

Based on the framework of ToPS, we design the first algorithm to output streams while satisfying local DP (LDP), which offers stronger privacy protection than DP. We call the resulting method ToPL. Under LDP, only the users know the true values and thus remove the dependence on the trusted central server. In ToPL, we use state-of-the-art LDP mechanisms for the Threshold optimizer and the Perturber. The design of ToPL relies on the findings in ToPS. In particular, ToPL implements the idea of minimizing both the bias and noise into the LDP primitive to find the optimal  $\theta$ . Moreover, we adapt existing LDP mechanisms to the problem tackled by this paper for better performance.

We implement both ToPS and ToPL, and evaluate them using four real-world datasets, including anonymized DNS queries, taxi trip records, click sreams and merchant transactions. We use the Mean Squared Error (MSE) over random range queries as the metric of performance evaluation. The experimental results demonstrate that our ToPS significantly outperforms the previous state-of-the-art algorithms. More specifically, the most significant improvement comes from our new method to find  $\theta$ . It contributes 4-8 orders of magnitude better utility than PAK. Even given the same reasonable  $\theta$ , ToPS can answer range queries  $100\times$  more accurately than PAK in terms of MSE. Putting the two together, ToPS improves 6-10 orders of magnitude in terms of utility over PAK.

**Contributions.** To summarize, the main contributions of this paper are threefold:

 We design ToPS for releasing real-time data streams under differential privacy. Its contributions include an EM-based algorithm to find the threshold, an on-line consistency algorithm, the use of a smoother to reduce the noise, and the ability to handle infinite streams.

- We extend ToPS to solve the problem in the more stringent setting of LDP and propose a new algorithm called ToPL. To the best of our knowledge, ToPL is the first attempt to publish streaming data under LDP.
- We evaluate ToPS and ToPL using several real-world datasets. The experimental results indicate that both can output stream pretty accurately in their settings, respectively. Moreover, ToPS outperforms the previous state-of-the-art algorithms by a factor of 6 – 10 orders.

**Roadmap.** In Section II, we present the problem definition and the background knowledge of DP and LDP. We present the existing solutions and our proposed method in Sections III. Experimental results are presented in IV. Finally we discuss related work in Section V and provide concluding remarks in Section VI.

## II. PROBLEM DEFINITION AND PRELIMINARIES

We consider the setting of publishing a stream of real values under differential privacy (DP). The length of the stream could be unbounded. And one wants to compute the aggregated estimates (e.g., prefix-sum and moving average) of the values. For instance, in [33], the authors gave an example of estimating average commute time in real-time where the stream consists of commute times of trips made by passengers.

## A. Formal Problem Definition

There is a sequence of readings  $V = \langle v_1, v_2, \ldots \rangle$ , each being a real number in the range of [0, B]. The goal is to publish a private sequence  $\tilde{V}$  of the same size as V in a way that satisfies differential privacy.

We measure the similarity between the private stream and the estimated stream using range queries, similar to [33]. Specifically, let V(i,j) be the sum of the stream from index i to j, i.e.,  $V(i,j) = \sum_{k=i}^{j} v_k$ . We want a mechanism that minimizes the expected squared error of any randomly sampled range queries

$$\mathbb{E}\left[\left(\tilde{V}(i,j) - V(i,j)\right)^2\right]. \tag{1}$$

## B. Differential Privacy

We follow the setting of [33] and focus on approaches that provide *event-level* DP [19].

**Definition 1** (Event-level  $(\epsilon, \delta)$ -DP). An algorithm  $\mathbf{A}(\cdot)$  satisfies  $(\epsilon, \delta)$ -differential privacy  $((\epsilon, \delta)$ -DP), where  $\epsilon, \delta \geq 0$ , if and only if for any two neighboring sequences V and V' and for any possible output set O,

$$\Pr\left[\mathbf{A}(V) \in O\right] \le e^{\epsilon} \Pr\left[\mathbf{A}(V') \in O\right] + \delta$$

where two sequences  $V = \langle v_1, v_2, \ldots, \rangle$  and  $V' = \langle v'_1, v'_2, \ldots, \rangle$  are neighbors, denoted by  $V \simeq V'$ , when  $v_i = v'_i$  for all i except one index.

As the algorithm **A** has access to the true sequence V, this model is called the centralized DP model (DP for short). For brevity, we use  $(\epsilon, \delta)$ -DP to denote Definition 1. When  $\delta = 0$ ,

which is the case we consider in this paper, we omit the  $\delta$  part and write  $\epsilon$ -DP instead of  $(\epsilon, 0)$ -DP.

Event-level DP means that information about each individual event remains private. In many cases, event-level DP is a suitable guarantee of privacy, e.g., individuals might be happy to disclose their routine trip to work while unwilling to share the occasional detour. The definition of event-level DP is also adopted by other works such as [11], [19], [9].

Compared to the centralized setting, the local version of DP offers a stronger level of protection, because each value is reported to the server in a perturbed form. The user's privacy is still protected even if the server is malicious. For each value  $\boldsymbol{v}$  in the stream of  $\boldsymbol{V}$ , we have the following guarantee:

**Definition 2**  $((\epsilon, \delta)$ -LDP). An algorithm  $\mathbf{A}(\cdot)$  satisfies  $(\epsilon, \delta)$ -local differential privacy  $((\epsilon, \delta)$ -LDP), where  $\epsilon, \delta \geq 0$ , if and only if for any pair of input values v, v', and any set O of possible outputs of  $\mathbf{A}$ , we have

$$\Pr[\mathbf{A}(v) \in O] \le e^{\epsilon} \Pr[\mathbf{A}(v') \in O] + \delta$$

Typically,  $\delta=0$  in LDP. Thus we simplify the notation and call it  $\epsilon$ -LDP. The notion of LDP differs from DP in that each user perturbs the data before sending it out and thus do not need to trust the server under LDP. In this paper, we focus on both the DP and the LDP setting.

## C. Mechanisms of Differential Privacy

We first review several primitives proposed for satisfying DP. We defer the descriptions of LDP primitives to Appendix C as our LDP method mostly uses the LDP primitives as blackboxes.

**Laplace Mechanism.** The Laplace mechanism computes a function f on the input V in a differentially private way, by adding to f(V) a random noise. The magnitude of the noise depends on  $\mathsf{GS}_f$ , the *global sensitivity* or the  $L_1$  sensitivity of f, defined as,

$$\mathsf{GS}_f = \max_{V \sim V'} ||f(V) - f(V')||_1$$

When f outputs a single element, such a mechanism  $\mathbf{A}$  is given below:

$$\mathbf{A}_f(V) = f(V) + \mathsf{Lap}\left(rac{\mathsf{GS}_f}{\epsilon}
ight)$$

In the definition above, Lap  $(\beta)$  denotes a random variable sampled from the Laplace distribution with scale parameter  $\beta$  such that  $\Pr[\mathsf{Lap}(\beta) = x] = \frac{1}{2\beta} e^{-|x|/\beta}$ . When f outputs a vector,  $\mathbf A$  adds independent samples of  $\mathsf{Lap}\left(\frac{\mathsf{GS}_f}{\epsilon}\right)$  to each element of the vector.

**Exponential Mechanism.** The exponential mechanism (EM) [30] samples from the set of all possible answers according to an exponential distribution, with answers that are "more accurate" being sampled with higher probability. This approach requires the specification of a quality functions q that takes as input the data V, a possible output o, and outputs a real-numbered quality score. The global sensitivity of the quality functions  $\mathsf{GS}_q$  is defined as:

$$\mathsf{GS}_q = \max_{o} \max_{V \simeq V'} |q(V, o) - q_i(V', o)| \tag{2}$$

The following method **A** satisfies  $\epsilon$ -differential privacy:

$$\Pr\left[\mathbf{A}_{q}(V) = o\right] = \frac{\exp\left(\frac{\epsilon}{2\operatorname{\mathsf{GS}}_{q}}q(V, o)\right)}{\sum_{o}' \exp\left(\frac{\epsilon}{2\operatorname{\mathsf{GS}}_{q}}q(V, o')\right)} \tag{3}$$

As shown in [30], if the quality function satisfies the condition that when the input dataset is changed from V to V', the quality scores of all outcomes change in the same direction, i.e., for any neighboring V and V'

$$(\exists_o q(V, o) < q(V', o)) \Longrightarrow (\forall_o' q(V, o') \le q(V', o')).$$

Then one can remove the factor of 1/2 in the exponent of Equation (3) and return i with probability proportional to  $\exp\left(\frac{\epsilon}{\mathsf{GS}_q}q(V,o)\right)$ . This improves the accuracy of the result.

## D. Composition Properties

The following composition properties hold for both DP and LDP algorithms, each commonly used for building complex differentially private algorithms from simpler subroutines.

**Sequential Composition.** Differential privacy is composable in the sense that combining multiple subroutines that satisfy differential privacy for  $\epsilon_1, \cdots, \epsilon_m$  results in a mechanism that satisfies  $\epsilon$ -differential privacy for  $\epsilon = \sum_i \epsilon_i$ . Because of this, we refer to  $\epsilon$  as the privacy budget of a privacy-preserving data analysis task. When a task involves multiple steps, each step uses a portion of  $\epsilon$  so that the sum of these portions is no more than  $\epsilon$ .

**Parallel Composition.** Given m algorithms working on disjoint subsets of the dataset, each satisfying DP for  $\epsilon_1, \dots, \epsilon_m$ , the result satisfies  $\epsilon$ -differential privacy for  $\epsilon = \max_i \epsilon_i$ .

**Post-processing.** Given an  $\epsilon$ -DP algorithm **A**, releasing  $g(\mathbf{A}(V))$  for any g still satisfies  $\epsilon$ -DP. That is, post-processing an output of a differentially private algorithm does not incur any additional loss of privacy.

The composition properties allow us to execute multiple differentially private computations and reason about the cumulative privacy risk. In our applications, we want to bound the total risk so we impose a total "privacy budget"  $\epsilon$  and allocate a portion of  $\epsilon$  to each private computation.

## III. PUBLISHING STREAMING DATA VIA DIFFERENTIAL PRIVACY

## A. Existing Work: PAK

For privately release binary stream, Dwork et al. [19] proposed an approach which is to inject Laplace noise to each data point in the stream and use hierarchical method for answering range queries. Perrier et al. [33] studied a more general setting where data points are sampled from [0, B], where B is some public upper bound. They observed that most of the values are concentrated below a threshold and proposed a method, which we call PAK for short (the initials of the authors' last names), to find such threshold and truncate data points below it reduce the scale of the injected Laplace noises. In PAK, the first m values are used to estimate a reasonable

threshold  $\theta$  which is smaller than the values' upper bound B under differential privacy. After that, the hierarchical method is used for estimating the stream statistics with the remaining n-m values (PAK assumes there are n observations).

**Phase 1: Finding the Threshold.** The intuition is that in many real world streaming data, the public upper bound B can be quite large, but most data points in the stream are much smaller than this bound. For instance, the largest possible purchase price of supermarket transactions is much larger than what a ordinary customer usually spends. After obtaining  $\theta$ , the following values in the stream are truncated to be no larger than  $\theta$ . Reducing the upper bound from B to  $\theta$  reduces the DP noise (via reducing sensitivity).

The work of PAK proposed a specially designed algorithm based on smooth sensitivity [32] to get the p-quantile (or p-percentile) as  $\theta$ , i.e., p% of the values are smaller than  $\theta$ . Here we describe the high-level idea of the algorithm.

As shown in [32], smooth sensitivity can be used to compute the median in a differentially private manner. It can be easily extended to privately release the p-quantile. In [33], the authors claim that the result of smooth sensitivity is unbiased, but they want the result to be larger than the real p-quantile. This is because if the estimated percentile is smaller, the truncation in the next phase will introduce more bias. As a result, PAK modifies the original smooth sensitivity method to guarantee that it is unlikely the result is smaller than the real p-quantile. As the details of the method is not directly used in the rest of the paper, we defer the details of both smooth sensitivity and the algorithm itself to Appendix A and B.

There are two drawbacks of this method. First, it requires a p value to be available beforehand. But a good choice of p actually depends on the dataset,  $\epsilon$ , and m. PAK simply uses p=99.5. As shown in our experiment in Section IV-D, it does not perform well in every scenario. Second, PAK bounds the probability that  $\theta$  is smaller than the p-quantile to be small. To ensure this, a positive bias is introduced to  $\theta$ , making  $\theta$  even larger.

**Phase 2: The Hierarchical Algorithm.** To achieve higher accuracy for range sums computed from the released private stream, the most straightforward way is to add independent noise generated through the Laplace distribution scaled to  $\theta$ , where  $\theta$  is the threshold obtained from the previous phase. However, this results in cumulative error (absolute difference from the true sum) of  $O(\theta\sqrt{n})$  after n observations.

To get rid of the linear dependency on  $\sqrt{n}$ , the hierarchical method is proposed [19], [9]. Given a string of length n, the algorithm first constructs a tree: the leaves are labelled by the intervals  $[1,1],[2,2],\ldots,[n,n]$  and each parent node is the union of intervals of its child nodes. To output the noisy count  $\tilde{V}(i,j)$ , the algorithm finds at most  $\log n$  nodes in the tree, whose union equal [i,j]. Given that the noise added to each node is only scaled to  $O(B\log n)$ , this results in an error of  $O(B\log^{1.5} n)$ .

As we will show in Section III-D, since the first proposal [19], there have been improvements to the hierarchical method [34], and it is straightforward to implement the new methods. Moreover, directly using the hierarchical method

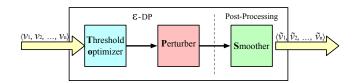


Fig. 1. Illustration of ToPS where differentially private noise is added to the incoming entry of a stream based on threshold  $\theta$  and final output is released after post-processing.

introduces overwhelming noise to the leaf nodes. We will demonstrate how to overcome this issue in Section III-E.

## B. Technical Overview of Our Approach: the Framework

The design of PAK was guided by asymptotic analysis. Unfortunately, for the parameters that are likely to occur in practice, the methods and parameters chosen by asymptotic analysis can be far from optimal, as such analysis ignores important constant factors. We advocate an approach that uses concrete error analysis to better guide the choice of methods and parameters.

In this section, we first deal with the threshold selection problem using the exponential mechanism (EM, as described in Section II-C), which satisfies pure DP (without introducing  $\delta$ ). Empirical experiments show its superiority especially in small  $\epsilon$  scenarios (which means compared to PAK, we can achieve the same performance with better privacy guarantees). We then introduce multiple improvements for the hierarchical methods including an online consistency method, and a method to reduce the noise in the lower levels of the hierarchy. We integrate all the components in a general framework, which stands for the combination of a Threshold optimizer, a Perturber, and a Smoother:

- Threshold optimizer: The threshold optimizer uses the input stream to find a threshold  $\theta$  for optimizing the errors due to noise and bias. It then truncates any incoming values by  $\theta$  and releases them to the perturber.
- Perturber: The perturber adds noise to each truncated incoming element of the stream, and releases noisy counts to the smoother.
- Smoother: The smoother performs further postprocessing on the noisy counts and output the final stream.

In the following subsections, we will instantiate the framework with differentially private algorithms. We call our method ToPS. Figure 1 gives the architecture of ToPS. It takes as input the raw stream  $V = \langle v_1, v_2, \ldots \rangle$  and outputs a private stream  $\tilde{V} = \langle \tilde{v}_1, \tilde{v}_2, \ldots \rangle$ . In this paper, we focus on the setting used in PAK. In particular, the threshold optimizer does not publishes the first m values, and instead uses them to obtain  $\theta$ . After the first m values, it sends  $\theta$  to the perturber and truncates any incoming value by  $\theta$ . We will discuss more about the flexibility of ToPS in the Section VI.

## C. Threshold Optimizer

We design a method, called threshold optimizer, to find the threshold  $\theta$  in a differentially private way. Given  $\theta$ , input

values greater than it will be truncated to  $\theta$ . However, this also introduces another error of bias due to the truncation. Thus choosing a good percentile that balances the two sources of errors is important.

PAK uses the approximate 99.5-th percentile as  $\theta$ . However, 99.5 is just a heuristic and may not work in every scenario. The choice should depend on factors such as  $\epsilon$  and others. For example, when the noise is small ( $\epsilon$  is relatively larger), we may want a larger percentile to reduce bias; but when  $\epsilon$  is smaller, we may want a small percentile.

We first formulate choosing  $\theta$  as an optimization problem that minimizes the combined error, and then introduce a function that approximates the error function while having a low sensitivity so that we can effectively use the Exponential Mechanism to solve it. We first consider the expected squared error of estimating a single value v. Assuming that  $\tilde{v}$  is the estimation of v, it is well known that the expected squared error is the summation of variance and the squared bias of  $\tilde{v}$ :

**Lemma 1.** 
$$\mathbb{E}\left[(\tilde{v}-v)^2\right] = \operatorname{Var}\left[\tilde{v}\right] + \operatorname{Bias}\left[\tilde{v}\right]^2$$
.

**Proof:** 

$$\begin{split} \mathbb{E}\left[\left(\tilde{v}-v\right)^{2}\right] = & \mathbb{E}\left[\left.\tilde{v}^{2}-2v\tilde{v}+v^{2}\right.\right] \\ = & \mathbb{E}\left[\left.\tilde{v}^{2}\right.\right]-2v\mathbb{E}\left[\left.\tilde{v}\right.\right]+v^{2} \\ = & \mathbb{E}\left[\left.\tilde{v}^{2}\right.\right]-\mathbb{E}\left[\left.\tilde{v}\right.\right]^{2}+\mathbb{E}\left[\left.\tilde{v}\right.\right]^{2}-2v\mathbb{E}\left[\left.\tilde{v}\right.\right]+v^{2} \\ = & \left(\mathbb{E}\left[\left.\tilde{v}^{2}\right.\right]-\mathbb{E}\left[\left.\tilde{v}\right.\right]^{2}\right)+\left(\mathbb{E}\left[\left.\tilde{v}\right.\right]-v\right)^{2} \\ = & \mathsf{Var}\left[\left.\tilde{v}\right.\right]+\mathsf{Bias}\left[\left.\tilde{v}\right.\right]^{2} \end{split}$$

Given a threshold  $\theta$  and privacy budget  $\epsilon$ ,  $\text{Var}\left[\tilde{v}\right] = \frac{2\theta^2}{\epsilon^2}$  and  $\text{Bias}\left[\tilde{v}\right] = \max(v - \theta, 0)$ .

Since we are using the hierarchical method [34] to publish data stream for answering range queries, we use the error estimations of the hierarchical method to instantiate Lemma 1. Qardaji et al. [34] show that there are approximately  $(b-1)\log_b(r)$  nodes to be estimated given any random query of range smaller than r, where b is the fan-out factor of the hierarchical method; and the variance of each node is  $\frac{2\theta^2}{\epsilon^2}$ . For bias, within range limitation r, a random query will cover around r/3 leaf nodes on average. Thus the combined error of the hierarchical method for answering random range query would be:

$$(b-1)\log_b^3(r)\frac{2\theta^2}{\epsilon^2} + \left(\frac{r}{3}\sum_{\theta < t < B} (\Pr[v=t](t-\theta))\right)^2$$
 (4)

For the exponential mechanism (EM) to be effective, the quality function should have a low sensitivity. However, if we directly use Equation 4 as the quality function, the sensitivity is large: a change of value from 0 to B will result in the increase of Equation 4 by  $(B-\theta)^2/9$ . Thus we choose to approximate the mean squared error in our quality function in the following ways.

Denote m as the number of values to be used in EM, and  $m_{\theta}$  as the number of values that are smaller than  $\theta$  from these m values, i.e.,

$$m_{\theta} = |\{i \mid v_i \le \theta, i \in [m]\}|,$$

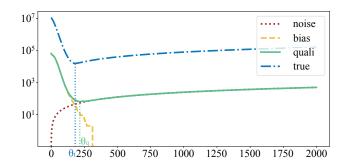


Fig. 2. Empirical comparison of the quality scores (Equation (5)) and the real measured errors and their minimum points  $\theta_q$  and  $\theta_t$  on a real-world datasets (DNS). We use  $\epsilon=0.1, m=2^{16}, r=2^{20}$ . The x-axes is the possible range upper bound, and the y-axes is the quality score or the measured error.

where  $[x] = \{1, 2, \ldots, x\}$  and |X| denotes the cardinality of set X. The first approximation method we use is to replace the variance and squared bias with standard deviation and bias. Second, we use  $c \cdot m_{\theta}$  to approximate  $\Pr[v = ](-\theta)$ . Here c is a constant that we will describe later. Finally, we multiply both the standard deviation and bias errors by  $-\frac{3m}{c \cdot r}$  to ensure the sensitivity is 1 (the sensitivity is defined in Equation (2)), and the quality score of the target is the highest. Equation 5 gives us the quality function:

$$q(V,\theta) = -\frac{3m}{c \cdot r} \sqrt{(b-1)\log_b^3(r) \frac{2\theta^2}{\epsilon^2}} - m_\theta$$
$$= -\frac{3m\theta}{cr\epsilon} \sqrt{2(b-1)\log_b^3(r)} - m_\theta. \tag{5}$$

The first term is a constant depending on  $\theta$  and is independent on the private data, while the second term has sensitivity of 1. The set of possible  $\theta$  values considered in the quality function  $q(V,\theta)$  should be a discrete set that covers the range [0,B]. The granularity of the set plays a role. If it is too coarse-grained (e.g.,  $\theta \in \{0,B/3,2B/3,B\}$ ), the method is inaccurate, because we are not sampling enough from the range, and the desired value might be far from the sampled values. On the other hand, if it is too fine-grained, the EM algorithm will run slowly, but it does not influence the accuracy much because of the normalization in EM (i.e., by the denominator in Equation (3)). In the experiment, we use all integers in the range of  $[B] = \{1,2,\ldots,B\}$  as the possible set of  $\theta$ .

One unexplained parameter in the quality function (5) is c. We choose c=60. The rationale for choosing c is that (a) We use  $m_{\theta}$  to approximate the bias term  $\Pr[v=t] (t-\theta)$ . This is an under-estimation, so we want to scale it up. (b) As we will describe later in Section III-D and III-E, the actual variance will be reduced by our newly proposed method (but there is no closed form for it); thus the variance formula overestimates the real variance, and we want to scale up the bias formula further to compensate for that. In practice demonstrated in the evaluation, we observe using  $3m_{\theta}$  to approximate  $\Pr[v=t] (t-\theta)$  is reasonable; moreover, the empirical improvement of our new method is roughly  $20 \times$ . Thus we set c=60.

Figure 2 illustrates the distribution of quality scores and measured errors on a dataset which is used in the experiments in Section IV. The dataset is a network streaming dataset, called DNS. We use  $\epsilon = 0.1, m = 2^{16}, r = 2^{20}$ , which

are the same as those parameters used in experiments. From Figure 2, we can see that the distributions between the true measured errors (Equation 4) and the corresponding quality scores (Equation 5) on two datasets are very close. The figure also illustrates the bias and variance factors of the quality scores. The two factors grow in opposite directions which makes a global minimum where the target  $\theta$  lies. In addition, we also show that the threshold  $\theta_q$  that minimizes our quality function Equation (5) is close to the target threshold  $\theta_t$  which minimizes the real measured errors (Equation 4). Therefore, the above empirical evaluation results show the capability of the threshold optimizer in finding accurate  $\theta$  values.

## D. Perturber

The perturber is implemented by the hierarchical method introduced in [19], which is the standard way of handling range queries. In this section, we investigate the hierarchical method and propose improvements for it.

1) Improving the Hierarchical Algorithm: We first present two straight-forward extensions of the hierarchical method.

**Better Fan-out.** First of all, we note that according to [34], using a fan-out b=16 instead of 2 (used in PAK) in the hierarchy can give better empirical result. This is done by optimizing the expected error numerically (instead of asymptotically). We thus use fan-out b=16 by default.

Handling Infinite Streams. Second, we extend the hierarchical method to support infinite streams. PAK requires a fixed length n in the beginning. After n observations, the algorithm will stop. To handle this infinite stream model, we propose to keep applying hierarchies. There is a question of how high the hierarchy should be. We can have higher hierarchies for handling larger queries. This is also what [9] proposed. But we note that most of the queries focus on small range. We thus propose to have an upper bound on the query range denoted by r and then split  $\epsilon$  equally to the  $h = \lceil \log_h r \rceil$  layers. The value of r stands for a limit below which most queries' range fall, and is determined externally. For each chunk of r observations, we output a height-h hierarchy. We note that each hierarchy handles a disjoint sub-stream of observations and thus this extension does not consume additional privacy budget because of the parallel composition property of differential privacy.

2) On-line Consistency: Given a noisy hierarchy, Hay et al. [24] proposed an efficient algorithm for enforcing consistency among the values in it. By enforcing consistency, the accuracy of the hierarchy can be improved. Note that this is a post-processing step and does not consume any privacy budget. Unfortunately, the algorithm is off-line and requires the whole hierarchy data to be saved (i.e., after every r observations). In this subsection, we propose an on-line algorithm to enforce consistency. That is, we can output the noisy streams promptly, and the noisy stream is consistent itself. Our method is built on [24]. So before going into details of our method, we first describe the consistency framework of Hay et al. [24].

**Off-line Consistency.** We use x to denote a node on the hierarchy H, and let  $\ell(x)$  to be the height of x (the height of a leaf is 1; and root is of height h). We also denote  $\operatorname{prt}(x), \operatorname{chd}(x)$ , and  $\operatorname{sbl}(x)$  to denote the children, parent, and

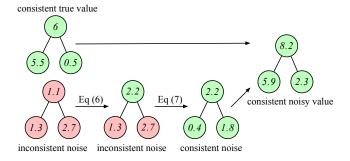


Fig. 3. Example of online consistency. The noisy value hierarchy breaks down into hierarchies of true value and noise where a parent node is supposed to be the sum of its child nodes. While true values are naturally consistent, this condition does not necessarily hold for the noise. The on-line consistency algorithm (application of Equation (7) and (6)) works only on the noise hierarchy to enforce the final output to be consistent. In the final output, only the leaf nodes are needed because the hierarchy is already consistent.

siblings of x, respectively. We use H(x) to denote the value corresponding to node x in the hierarchy. The first step updates the values in H from bottom to top. The leaf nodes remain the same; and for each height- $\ell$  node x, where  $\ell$  iterates from 2 up to h, we update it

$$H(x) \leftarrow \frac{b^l - b^{l-1}}{b^l - 1} H(x) + \frac{b^{l-1} - 1}{b^l - 1} \sum_{y \in \mathbf{chd}(x)} H(y).$$
 (6)

We then update the values again from top to bottom. This time the value on root remains the same, and for each height- $\ell$  node x, where  $\ell$  iterates from  $\ell$  down to 1, we update it

$$H(x) \leftarrow \frac{b-1}{b}H(x) + \frac{1}{b}\left(H(\operatorname{prt}(x)) - \sum_{y \in \operatorname{sbl}(x)} H(y)\right). \tag{7}$$

An On-line Algorithm. Our method is derived from the observation that the noisy estimates can be decomposed into true values and pure noise. The idea of our approach is to generate all the required noises in advance and then enforce consistency among the noises. During the on-line publishing, we can just add the consistent noises on the final results to ensure consistency.

More formally, we decompose the noisy values in H into two parts: the true values and the pure noises, denoted as T and N, respectively. N is the independent noise from the Laplace mechanism (described in Section II-C). T is defined by induction: for a leaf x, T(x) corresponds to one true value v, and for a node x in a higher level,  $T(x) = \sum_{y \in \operatorname{chd}(x)} T(y)$ . The true values from T are consistent naturally. Thus if N is consistent, H is also consistent. To see it, consider any internal node x,

$$\begin{split} &H(x) = T(x) + N(x) \\ &= \sum_{y \in \operatorname{chd}(x)} T(y) + \sum_{y \in \operatorname{chd}(x)} N(y) \\ &= \sum_{y \in \operatorname{chd}(x)} (T(y) + N(y)) = \sum_{y \in \operatorname{chd}(x)} H(y). \end{split}$$

In the on-line consistency algorithm, we internally generate the noise hierarchy N and run the steps given in Equations (6) and (7) to make N consistent. After this pre-processing step, we can ignore the higher-layers of the hierarchy, and use only the leaf nodes which add consistent noise to each individual value. This is because the results from the higher-layers are already consistent with those from the leaves, thus it suffices to only output the most fine-grained result. Figure 3 gives an example using a simple binary tree.

We prove that H(x) resulted from the online-algorithm is equivalent to that of the off-line algorithm.

**Theorem 1.** The on-line consistency algorithm gives the same result as the off-line consistency algorithm.

*Proof:* We first examine the bottom-up update step. According to Equation (6), the updated N(x) equals to

$$\frac{b^l - b^{l-1}}{b^l - 1} N(x) + \frac{b^{l-1} - 1}{b^l - 1} \sum_{y \in \operatorname{chd}(x)} N(y).$$

Adding T(x) to it, we have the updated H(x) equals to

$$\begin{split} &\frac{b^{l}-b^{l-1}}{b^{l}-1}N(x)+\frac{b^{l-1}-1}{b^{l}-1}\sum_{y\in\operatorname{chd}(x)}N(y)+T(x)\\ &=\frac{b^{l}-b^{l-1}}{b^{l}-1}(N(x)+T(x))+\frac{b^{l-1}-1}{b^{l}-1}\left(\sum_{y\in\operatorname{chd}(x)}N(y)+T(x)\right)\\ &=\frac{b^{l}-b^{l-1}}{b^{l}-1}(N(x)+T(x))+\frac{b^{l-1}-1}{b^{l}-1}\sum_{y\in\operatorname{chd}(x)}(N(y)+T(y)) \end{split}$$

$$=\frac{b^l-b^{l-1}}{b^l-1}H(x)+\frac{b^{l-1}-1}{b^l-1}\sum_{y\in \operatorname{chd}(x)}H(y).$$
 Equation (8) is because of the consistency of  $T$  so that

 $T(x) = \sum_{y \in \operatorname{chd}(x)} T(y)$ . This gives the same result as if we run the off-line consistency algorithm. Similarly, during the top-down update step (in Equation (7)), we have the updated N(x) equals to

$$\frac{b-1}{b}N(x) + \frac{1}{b}\left(N(\operatorname{prt}(x)) - \sum_{y \in \operatorname{\mathbf{sbl}}(x)} N(y)\right).$$

Adding T(x) to it, we have the updated H(x) equals to

$$\frac{b-1}{b}N(x) + \frac{1}{b}\left(N(\operatorname{prt}(x)) - \sum_{y \in \operatorname{sbl}(x)} N(y)\right) + T(x)$$

$$= \frac{b-1}{b}(N(x) + T(x)) + \frac{1}{b}\left(N(\operatorname{prt}(x)) - \sum_{y \in \operatorname{sbl}(x)} N(y) + T(x)\right)$$

$$= \frac{b-1}{b}H(x) + \frac{1}{b}\left(H(\operatorname{prt}(x)) - \sum_{y \in \operatorname{sbl}(x)} H(y)\right). \tag{9}$$
which is some result as if  $x \in \operatorname{spl}(x)$ 

which is same result as if we run the off-line consistency algorithm. Equation (9) holds is also because of the consistency of T so that  $T(x) = T(\operatorname{prt}(x)) - \sum_{y \in \operatorname{Sbl}(x)} T(y)$ .

## E. Smoother

In this section, we introduce a smoother to further improve the utility of the algorithm. The high-level idea of the smoother is to replace the individual estimate given by the perturber with predicted values. By doing this, the perturber can focus on only the aggregated summations and be more accurate.

More specifically, the original perturber works on h layers (although only the leaf nodes will be outputted, we draw Laplace noise from the other higher layers to make the noise consistent; and having a hierarchy helps when answering range queries). The smoother replaces some s lower levels of H with predicted values, and the perturber runs on the upper h-slayers. Next, we show how to determine s.

**Optimizing** s. Selecting s is important. A larger s results in smaller noise errors: because there are now h-s layers in the hierarchy, each layer will receive more privacy budget according to sequential composition (given in Section II-D). On the other hand, a larger s probably leads to a larger bias (because we are only doing the actual estimate once every  $b^s$ values; others in between are independent of, and thus can be far from, the true values). Choosing a good value of s thus is a balance between noise errors and bias. Unfortunately, the bias depends on the true data distribution. Thus we assume the bias for each value is approximately half of the truncation limit  $\theta$ . We use the following equation to approximate the squared error:

$$(b-1) (\log_b(r) - s)^3 \frac{2\theta^2}{\epsilon^2} + \frac{b^s}{2} \frac{\theta^2}{4}.$$
 (10)

Given the values of  $\epsilon$  and r, s can be computed by minimizing the above error.

**Smoothing Method.** Given s, we now describe choices of implementing the smoother. Denote  $u_1, u_2, \ldots$  as the noisy estimates given by the leaves of the perturber (or the (s+1)th levels of the original hierarchy). Each  $u_i$  is the noisy sum of  $b^s$  values. Let  $u_0 = \frac{1}{2}b^s\theta$  (initially there is no estimations from the hierarchy; we thus use half of the threshold as mean). The smoother will take the sequence of u and output the final result  $\tilde{v}_i$  for each input value. Let  $t = \lceil i/b^s \rceil$ , we consider several functions:

- 1) Recent smoother:  $\tilde{v}_i = u_t/b^s$ . It takes the mean of
- the most recent output from the perturber. Mean smoother:  $\tilde{v}_i = \frac{1}{b^s t} \sum_{j=0}^t u_j$ . It takes the mean of the output from the perturber up until the moment. 2)
- Median smoother:  $\tilde{v}_i = \text{median}(u_1, \dots, u_t)$ . Similar to the mean smoother, the median smoother takes the median of the output from the perturber up until the
- Moving average smoother:  $\tilde{v}_i = \frac{1}{b^s w} \sum_{j=t+1-w}^t u_j$ . Similar to the mean smoother, it takes the mean over 4) the most recent w outputs from the perturber. When t+1 < w, we use the average of the first t+1 values of u divided by  $b^s$  as  $\tilde{v_i}$ .
- Exponential smoother:  $\tilde{v}_i = \frac{u_0}{b^s}$  if t = 0, and  $\tilde{v}_i = \alpha \frac{u_t}{b^s} + (1 \alpha)\tilde{v}_{i-b^s}$  if t > 0, where  $0 \le \alpha \le 1$  is the smoothing parameter. The exponential smoother

put more weight on the more recent values from the hierarchy.

## F. Outputting Streams in the LDP Setting

In this section, we describe our method ToPL for outputting streaming data in the local DP (LDP) setting. To the best of our knowledge, this is the first algorithm that deals with this problem in LDP.

In LDP, users perturb their values locally before sending them to the server, and thus do not need to trust the server. Applying to the streaming values in our setting, each value should be perturbed before being sent to the server. What the server does is only post-processing of the perturbed reports. After that, the results can be used or shared with other parties.

Our method follows the design framework of ToPS. There is a threshold optimizer to find the threshold based on optimal estimated error; and the threshold is used to truncate the users' values in the later stage. Different from the centralized DP setting, in the local setting, the obtained threshold will be shared with the users so that they can truncate their values locally. The perturber section is also run within each user's local side, because of the privacy requirement that no other parties other than the users themselves can see the true data. There is no smoother section. In what follows, we describe the construction for the threshold optimizer and the perturber.

1) Design of the Threshold Optimizer: In LDP, each user only has a local view (i.e., they only know their own data; no one has a global view of the true distribution of all data), thus there is no Exponential Mechanism (EM) (described in Section II-C) that we can use as in the DP setting (the threshold optimizer in Section III-C). Instead, most existing LDP algorithms rely on frequency estimation, i.e., estimation of how many users possess each value, as what the Laplace mechanism does in DP. We also rely on the frequency estimation to find the optimal threshold. Although the distribution estimation is more informative, it is actually less accurate than EM because more noise needs to be added.

Frequency Estimation in LDP. Li et al. [29] propose the Square Wave mechanism (SW for short) for ordinal and numerical domains. It extends the idea of randomized response [44] in that values near the true value will be reported with high probability, and those far from it have a low probability of being reported. The server, after receiving the reports from users, runs a specially designed Expectation Maximization algorithm to find an estimated density distribution that maximizes the expectation of observing the output. For the completeness, we describe details about SW in Appendix C1.

Optimized Threshold with Estimated Distribution. To find the threshold, the baseline method is to find a specific percentile to find the a threshold  $\theta$ . This method is used for finding frequent itemset [42]. Based on the lessons learned from the threshold optimizer in the DP setting, we use the optimization equation given in Lemma 1 to find  $\theta$ .

Specifically, denote  $\tilde{f}$  as the estimated distribution where  $\tilde{f}_t$  is the estimated frequency of value t. Here the set of all possible t to be considered can no longer be  $[B] = \{1, 2, \dots, B\}$ . Instead, we sample 1024 values uniformly from [B]. This is

because SW uses the Expectation Maximization algorithm, and a large domain size makes it time and space consuming. Similar to Equation (4) considered in the DP setting, we use an error formula:

$$\frac{r}{3} \cdot \mathsf{Var}\left[\tilde{v}\right] + \frac{r^2}{24} \left( \sum_{\theta < t < B} \tilde{f}_t(t - \theta) \right)^2 \tag{11}$$

where  $\frac{r}{3}$  is the expected size of the range query. Here  $\text{Var}\left[\tilde{v}\right]$  denotes the variance of estimating v, which we will describe in the perturber. It is multiplied by  $\frac{r}{3}$  because in expectation, each random range query will involve  $\frac{r}{3}$  values, and each of them is estimated independently. For the second part of Equation (11), it can be calculated directly with SW (without any approximation, as in the case of Equation (5)). The multiplicative coefficient  $\frac{r^2}{24}$  is the averaged case over all possible range queries. That is, denote k as the range of a query, there are r-k+1 range-k queries within a limit r, and there are all together  $\sum_{k=1}^{r}(r-k+1)$  possible queries. For each of them, we have a  $k^2$  coefficient in the squared bias. Thus, we have  $\frac{\sum_{k=1}^{r}(r-k+1)k^2}{2r(r+1)} = \frac{(r+1)(2r+1)}{12} - \frac{r(r+1)}{8} \approx \frac{r^2}{24}$  as the average-case coefficient. Finally, we find a  $\theta$  that minimizes Equation (11).

Using SW as a White Box. One thing to note is that, SW is proposed for estimating smoothed distributions, while in our case, the distribution is very skewed, because the majority of values are expected to be concentrated below a threshold. To make SW output reasonable thresholds, instead of using SW as a black-box, we make the following modifications. First, we eliminate the smoothing step from SW, because we observe that in some cases, the smoothing operation will "push" the density to the two ends, which is unlikely to happen. Second, we add a post-processing step to prune the small densities outputted by SW. In particular, we find the first w so that  $\forall i < 5$ ,  $\tilde{f}_{w+i} < 0.1\%$ . This is a signal that the density after w will converge to 0. We thus replace the estimated density after w with 0. In the experiment, we observe that the two steps help finding a more accurate  $\theta$ .

2) Design of the (Local) Perturber: After obtaining the threshold  $\theta$ , the server sends  $\theta$  to all users. When a user reports a value, it will first be truncated. The user then reports the truncated value using the Hybrid mechanism. The method is described in Appendix C2. It can estimate v with worst-case variance of given in Equation (13), which can be plugged into Equation (11) to find  $\theta$ . Note that the reports are unbiased by themselves. So to answer a range query, we just need to sum up values from the corresponding range, and there is no need for a smoother.

## G. Discussion

We claim that ToPS and ToPL satisfy DP and LDP, respectively. For ToPS, the perturber uses  $\epsilon/h$  to add Laplace noise to each layer of the the hierarchical structure. By sequential composition, the overall data structure satisfies  $\epsilon$ -DP. To find the threshold, ToPS uses a disjoint set of m observations and runs an  $\epsilon$ -DP algorithm. Due to the parallel composition property of DP, the threshold optimizer and the perturber together satisfy  $\epsilon$ -DP. The on-line consistency algorithm and the smoother's operations are post-processing procedures and

do not affect the privacy guarantee. For ToPL, similarly, it satisfies  $\epsilon$ -LDP as we only use existing LDP primitives.

The design of our method is inspired by PAK. But we include some unique design choices to handle the problem. The difference between the two methods are as follows. PAK has two phases, the threshold finder and the hierarchical method. The biggest drawback of PAK is that it focuses on methods that provide good theoretical bound while ignored the empirical performance. We improve both phases using either new methods or methods from existing literature. As we use EM, our algorithm satisfies  $\epsilon$ -DP while PAK satisfies  $(\epsilon, \delta)$ -DP. Moreover, we introduce a smoother that further improves accuracy.

## IV. EXPERIMENTAL EVALUATION

The experiment includes four phases. First, we give a high-level end-to-end evaluation of the whole process. Second, we fix a truncation threshold and evaluate the performance of the hierarchical method. Third, we fix the hierarchical method and test different algorithms that gives the truncation threshold. Fourth, we evaluate the performance in the local setting.

#### A. Evaluation Setup

**Datasets.** Our experiments are conducted on four real world datasets.

- DNS: This dataset is extracted from a set of DNS query logs collected by a campus resolver with all user ids and source IP addresses removed. <sup>1</sup> This dataset includes 14 days of DNS queries. The number of queries is estimated, which can be used to assess Internet usage in a region.
- Fare [2]: New York city taxi travel fare. We use the Yellow Taxi Trip Records for January 2019.
- Kosarak [1]: A dataset of click streams on a Hungarian website that contains around one million users and 41270 categories. The data is formatted so that each user click through multiple categories. We take it as a streaming data and use the size of click categories as the value of the stream.
- POS [49]: A dataset containing merchant transactions of half a million users and 1657 categories. We use the size of the transaction as the value of the stream.

Table I gives the distribution statistics of the datasets.

TABLE I. DATASET CHARACTERISTICS

Dataset	n	max	$p_{100}$	$p_{85}$	$p_{95}$	$p_{99.5}$	avg
DNS	1141961	2000	617	63	85	135	37.9
Fare	8704495	30000	26770	440	1036	2037	279.9
Kosarak	990002	41270	2498	10	28	133	8.1
POS	515597	1657	165	13	21	39	7.5

**Metrics.** To evaluate the performance of different methods, we use the metric of Mean Squared Error (MSE) of answering

randomly generated queries. In particular, we measure

$$MSE(Q) = \frac{1}{|Q|} \sum_{(i,j) \in Q} \left| \tilde{V}(i,j) - V(i,j) \right|^2.$$
 (12)

where Q is the set of the randomly generated queries. It reflects the analytical utility measured by Equation 1 from Section II-A. We set  $r=2^{20}$  as the maximal range of any query.

**Methodology.** The prototype was implemented using Python 3.7.3 and numpy 1.15.3 libraries. The experiments were conducted on servers running Linux kernel version 5.0 with Intel Xeon E7-8867 v3 CPU @ 2.50GHz and 576GB memory. For each dataset and each method, we randomly choose 200 range queries and calculate their MSE. We repeat each experiment 100 times and report the result of mean and standard deviation. Note that the values of standard deviation are typically very small, and barely noticeable in the figures.

## B. Results of the Whole Process

First of all, as a case study, we visualize the estimated stream of our method ToPS on the DNS dataset. Figure 4 shows the result. The top row shows the performance of ToPS. As a comparison, we also plot the visualization of PAK in the bottom row. We run algorithms only once for each setting to demonstrate the real-world usage. Similar to the setting of PAK, in ToPS, we use the first m=65,536 observations to obtain the threshold  $\theta$  (we will show later that ToPS does not need this large m observations to be blocked). As can be observed from Figure 4, our method ToPS can give fairly accurate predictions when  $\epsilon$  is pretty small. On the other hand, PAK, even though under a larger  $\epsilon$ , still performs worse than ToPS. Note that to obtain the threshold  $\theta$ , PAK satisfies  $(\epsilon, \delta)$ -DP. Our ToPS satisfies pure  $\epsilon$ -DP during the whole process.

We then compare the performance of ToPS and PAK with our metric of MSE given in Equation (12), and show the results in Figure 5. Between ToPS and PAK, we also include two intermediate methods that replace the Phase 1 (finding  $\theta$ ) and 2 (hierarchical method) of PAK by our proposed method EM-E (used in threshold optimizer) and  $\hat{H}^c_{16}$  (used for the perturber and smoother together), respectively, to demonstrate the performance boost due to our new design (we will evaluate the two phases in more details in later subsections). From the figure, we can see that the performance of all the algorithms gets better as  $\epsilon$  increases, which is as expected. Second, our proposed ToPS can outperform PAK by 7 to 11 orders of magnitude. Third, the effect (in terms of improving utility) using EM-E is much more significant than using  $\hat{H}_{16}^c$ . Interestingly, the performance of ToPS and EM-E is similar in the Fare and Kosarak datasets. This is because in these cases, the bias (due to truncation by  $\theta$ ) is dominant.

## C. Performance of Publishing the Stream

Several components contribute to the promising performance of ToPS. To demonstrate the precise effect of each of them, we next analyze them one by one in the reverse order. Namely, we first fix other things and compare different designs of the smoother and the perturber. Then, we analyze the methods of obtaining the threshold  $\theta$  in Section IV-D.

<sup>&</sup>lt;sup>1</sup>The data collection process has been approved by the IRB of the campus.

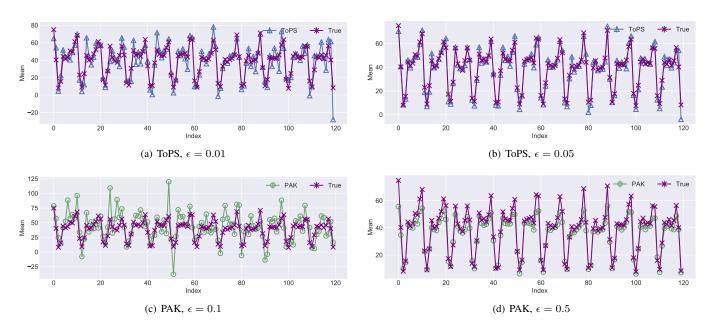


Fig. 4. Visualizations of the DNS stream. The x-axes corresponds to the time, and the y-axes denotes the moving average. Our ToPS at  $\epsilon = 0.01$  can output predictions that are pretty close to the ground truth. PAK gives more noisy result even with larger  $\epsilon$  values.

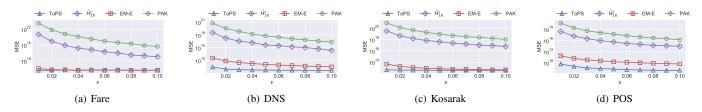


Fig. 5. Comparison between PAK and ToPS when answering range queries. We also include two intermediate methods EM-E (our proposed threshold optimizer) and  $\hat{H}_{16}^c$  (our proposed the perturber and smoother) that replace the corresponding two phases of PAK to demonstrate the performance boost due to our new design.

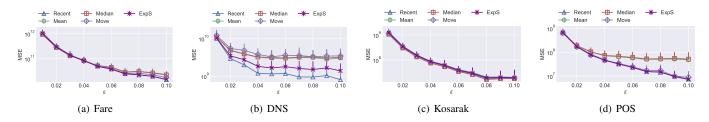


Fig. 6. Evaluation of different smoothing techniques. We vary  $\epsilon$  from 0.01 to 0.1 in the x-axis. The y-axis shows the query accuracy (MSE).

To make the comparison clear, we set  $\theta$  to be the 95-th percentile of the values. Moreover, we assume the true values are no larger than  $\theta$  (the ground truth is truncated). We will compare the performance of different methods in obtaining  $\theta$  in Section IV-D.

Comparison of Different Smoothers. Fixing a threshold  $\theta$  and the the hierarchical method optimized in Section III-D, we now compare the performance of five smoother algorithms listed in Section III-E (note that the smoothers will replace the  $16^s$  values where s is given in Equation (10)). In Figure 6, we vary the value of  $\epsilon$  from 0.01 to 0.1 and plot the MSE of the smoothers. Overall the performance gets better when  $\epsilon$  increases, which is as expected. The difference of them are

very small in the Fare and Kosarak datasets. This is because there is no clear pattern in these datasets. In the DNS dataset, Recent performs better than others, as the data is stable in the short term. In POS, the method of Mean and Median performs worse than the other three. This is because Mean and Median consider all the history (all the previous  $u_i$  values given from the hierarchy, as described in Section III-E), while the other methods considers more recent results. And due to the stability property in the dataset (similar to the case of DNS), methods that utilize the recent output (i.e., the more recent  $u_i$ ) will perform better. As Recent performs the best in DNS, and is among the best in other datasets, we use it as the default

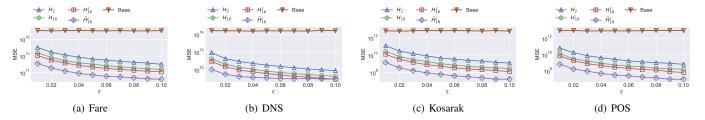


Fig. 7. Evaluation of different methods of outputting the stream. We vary  $\epsilon$  from 0.01 to 0.1 in the x-axis and plot the query accuracy (MSE) in the y-axis.

smoother algorithm.

**Comparison of Different Hierarchical Algorithms.** To demonstrate the precise effect of each design detail, we line up several intermediate protocols for the hierarchy and threshold, respectively. For the hierarchy component, we evaluate:

- $H_2$ : Original binary tree used in PAK.
- $H_{16}$ : A good fan-out b = 16 is used in the hierarchy.
- $H_{16}^c$ : The optimal  $H_{16}$  with consistency method.
- $\hat{H}_{16}^c$ : We use the hat notation to denote the Recent smoother. It is built on top of  $H_{16}^c$ .
- Base: A baseline method that always output 0. Base is used to understand whether each method gives meaningful result.

Figure 7 gives the result varying  $\epsilon$  from 0.01 to 0.1. First of all, we see that all methods (except the baseline) give better accuracy as  $\epsilon$  increases, which is as expected. Moreover, except  $H_{16}^c$ , the performance of all other methods increases by a factor of  $100\times$  when  $\epsilon$  increases from 0.01 to 0.1. This is also indicated by the analysis (variance is proportional to  $1/\epsilon^2$ ). Comparing each method, we observe that using the optimal branching factor ( $H_{16}$  versus  $H_2$ ) in this case also gives a  $5\times$ improvement, and by adopting the consistency step, we have another roughly  $2 \times (H_{16}^c \text{ versus } H_{16})$  of accuracy boost. For  $\hat{H}_{16}^c$ , we see a constant  $10\times$  improvement over  $H_{16}^c$  except in the DNS dataset, where  $H_{16}^c$  performs roughly the same as  $H_{16}^c$  when  $\epsilon > 0.08$ . The reason here is that the error of  $\hat{H}_{16}^c$ is composed of two parts, the noise error from the constraint of DP, and the bias error of outputting the predicted values. When  $\epsilon$  is large, the second error dominates the first one in the DNS dataset.

## D. Comparison of Threshold Phase

Having examined the performance of different methods of outputting the stream, we now switch gear to look at the algorithms that are used to find the threshold.

**Setup.** Following [33], we use the first m values to obtain the threshold. To eliminate the unexpected influence of distribution change, for now, we use the same m values to build the hierarchy using  $\hat{H}_{16}^c$  (with the best smoothing strategy called recent smoother).

**No Single Quantile Works Perfectly for All Scenarios.** We first want to show that there is no single *p*-quantile that can work perfectly for every scenario. We use the true *p*-quantile

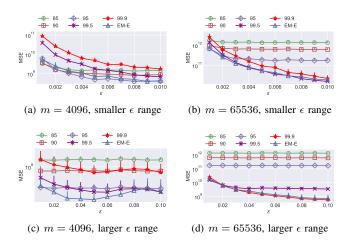


Fig. 8. MSE of answering range queries using  $\hat{H}_{16}^c$  on DNS dataset. The true p-th percentile for different p values are evaluated. We also include EM-E, which uses  $\epsilon=0.05$  to obtain  $\theta$ .

for  $p \in \{85, 90, 95, 99.5, 99.9\}$  and test in different scenarios (with different  $\epsilon$  and m values). We also include our threshold optimizer (EM-E) which is introduced to find a threshold only based on the estimated error and set  $\epsilon = 0.05$  for it. Figure 8 shows the results of answering range queries given these true percentiles. Several findings are demonstrated from the figures. First of all, the performance gets better when  $\epsilon$  gets larger for all p values. Second, in some cases, the performance improvement is negligible with respect to  $\epsilon$  (e.g., p=80 and 85 in Figure 8(b) and p=85, 90 and 95 Figure 8(d)). This is because in these scenarios, p is set too small, making the bias dominates the noise error. Third, our threshold optimizer with  $\epsilon=0.05$  can achieve a similar performance with the optimal p-quantile.

**Varying**  $\epsilon$ . We then compare different methods of outputting  $\theta$ . We evaluate:

- EM-E: The threshold optimizer in ToPS. It does not require a percentile.
- S-PAK: The smooth sensitivity method used in PAK.
   We use p = 99.5, as used in [33].
- S-P: The original smooth sensitivity method. Similar to S-PAK, we also use p = 99.5.

In Figure 9, we compare with existing differentially private methods on finding the threshold  $\theta$ . We vary the value of  $\epsilon$  for obtaining  $\theta$ . The result is the MSE of answering range queries using  $\hat{H}_{16}^c$ . Note that to make the comparison more clear, we

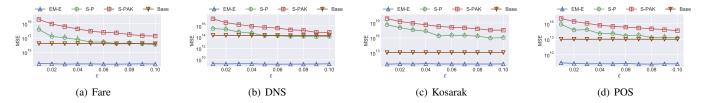


Fig. 9. Evaluation of different methods to find the threshold  $\theta$ . We vary  $\epsilon$  from 0.01 to 0.1 in the x-axis. The y-axis shows the MSE of answering range queries using  $\hat{H}_{6}^{r}$  (to make comparison clear, we use a fixed  $\epsilon = 0.05$  for it). Base is a baseline method that always output 0.

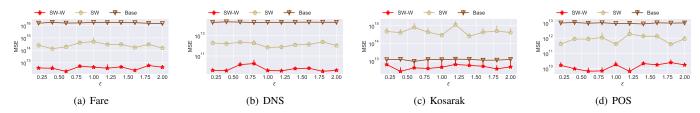


Fig. 10. LDP evaluation of different methods of outputting the threshold. We vary  $\epsilon$  from 0.2 to 2 in the x-axis. The y-axis shows the query accuracy (MSE).

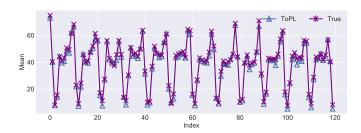


Fig. 11. Visualizations of the DNS stream. The x-axes correspond to the time, and the y-axes denote the moving average. Our ToPL at  $\epsilon=1$  can output predictions that are pretty close to the ground truth.

fix  $\epsilon=0.05$  in  $\hat{H}^c_{16}$ . In all the datasets, our proposed EM-E performs much better than existing methods. Moreover, the performance does not change much when  $\epsilon$  increases. This is because the output of EM-E is stable even in small  $\epsilon$  domains. Finally, both S-PAK and S-P perform worse than the baseline method, which always give 0 on any input values, indicating the  $\theta$  given by them is too large.

## E. Performance of ToPL

In this section, we evaluate the LDP algorithm ToPL. We first check the methods to find  $\theta$ . Following the setting of ToPS, we also use the first m=65,536 observations to obtain the threshold  $\theta$ . The difference from the DP setting is that due to the higher amount of noise of LDP, we vary  $\epsilon$  in a larger range (from 0.2 to 2). We remind the readers that our method find  $\theta$  by using the square wave (SW) mechanism to estimate the distribution, and then minimizing Equation (11). Moreover, our approach modifies SW to exploit the prior knowledge that the distribution is skewed. We use SW-W to denote this method. In addition, we also include another baseline, which uses SW as a black box (and we use SW to denote it).

Figure 10 shows the performance of finding  $\theta$  in ToPL. We vary the  $\epsilon$  used to find  $\theta$  from 0.2 to 2 while fixing  $\epsilon$  used to output the stream at 1. Similar to the DP setting, we did not observe the performance improvement because the  $\epsilon$ 

for the perturber is fixed, and both SW and SW-W can output pretty stable  $\theta$  with different  $\epsilon$  values. Comparing SW and SW-W, we can observe that our method SW-W can outperform SW as well as the baseline in all the four datasets, and the improvement (measured by MSE of answering random range queries) is around 2 orders of magnitude.

In Figure 11, we visualize the estimated stream of our method ToPL on the DNS dataset using  $\epsilon=1$ . We run algorithms only once to demonstrate the real-world usage. As can be seen from this figure, our method ToPL can give pretty accurate predictions.

## V. RELATED WORK

## A. Dealing with Streaming Data in Differential Privacy

A number of solutions have been proposed for solving the problem of releasing real-time aggregated statistics under differential privacy (DP). Here, besides the PAK approach [33], we briefly describe several other related works.

**Event-level DP.** The first line of work is the hierarchical method for handling binary stream. Event-level DP is satisfied in this case. Dwork et al. [19] proposed a differentially private continual counter release algorithm over a fixed length binary stream with a bounded error  $O\left(\left(\log^{1.5} n\right)/\epsilon\right)$  at each time step, where n is the stream length. Chan et al. [9] studied the same problem and proposed a release algorithm based on binary tree which achieves a similar error bound without requiring a prior knowledge of the stream length. There is also on on-line consistency method proposed in [9]. But that method focuses on enforcing integrality of the output (because it handles the binary setting), while our method minimizes the overall noise error.

Chen et al. [11] also considers the event-level DP, but in a different setting. In particular, the data can be any number instead of only a binary number. The data model is the same as what our paper and [33] consider. Moreover, [11] considers the application of outputting the number of connections to a wifi hotspot, and protect any event of a single user (i.e., whether a user connects to the hotspot). Thus, the sensitivity is only 1 (while in our setting, the original sensitivity is the maximal possible value B, and we first find a threshold  $\theta < B$ ). The method in [11] partitions the stream into a series of intervals so that values inside each interval are "stable", and publish the median of each interval.

User-level DP. There is also work that focuses on providing the notion of user-level DP. Because the user-level DP is more challenging, proposals under this setting relies more on the auxiliary information. In particular, Fan et al. [22] proposed to release perturbed statistics at sampled timestamps and uses the Kalman filter to predict the non-sampled values and correct the noisy sampled values. It takes advantage of the seasonal patterns of the underlying data. Another direction is the offline setting, where the server has the global view of all the values first, and then releases a streaming model satisfying DP. In this setting, [4] proposes an algorithm based on Discrete Fourier Transform (DFT), and [36] further incorporate sampling, clustering, and and smoothing into the process.

w-event-level DP. To balance the privacy loss and utility loss between user-level and event-level privacy models, relaxed privacy notions are proposed. Bolot et al. [8] extended the binary tree mechanism on releasing perturbed answers on sliding window sum queries over infinite binary streams with a fixed window size and using a more relaxed privacy concept called decayed privacy. Kellaris et al. [26] proposes w-event DP and two new privacy budget allocation schemes (i.e., split  $\epsilon$  into individual events) to achieve it. More recently, [39] works explicitly for spatiotemporal traces, and improves Kellaris et al. [26] by adaptively allocating privacy budget based on the similarity of the data sources; and [23] improves Kellaris et al. [26] using a similar approach of Fan et al [22]: sampling representative data points to add noise, and smoothing to reconstruct the other data points.

In our work, we follow the event-level privacy model and dealing with a more extended setting where data points are from a bounded range instead of the binary domain, and publish the stream in an on-line manner. Moreover, we do not rely on any pattern to exist in the data; and we propose methods for both DP and LDP.

## B. Dealing with Streams in Local DP

In the local DP setting, most existing work focus on estimating frequencies of values in the categorical domain [21], [7], [6], [40], [46], [3]. These techniques can be applied to other applications such as heavy hitter identification [6], [43], [37] frequent itemset mining [35], [42], multi-dimension data estimation [41], [48], [13], [14]. In the numerical/ordinal setting, previous work [16], [38] mostly focused on estimating mean. Recently, Li et al. [29] proposed the square wave mechanism for the more general task of estimating the density.

To the best of our knowledge, there are two methods that deals with streaming data in user-level LDP; and they work in different data models. In particular, [20] assumes the users' values can only change by 1, and there are at most k changes in a given time range. [25] assumes all users' values are repeatedly drawn from some distributions, e.g., a bit from a Bernoulli distribution.

## C. Other Related Topics

Choosing the Threshold. Many DP problems require choosing a threshold for truncation. Examples include: (1) Choosing contribution limit in transactional data such as finding frequent itemset [47], [28]. (2) Choosing configuration for publishing graph degree sequence under node-DP (e.g., sensitivity in social network is typically very large) [15]. (3) Answering aggregation queries over relational database (each user's contribution, a.k.a., sensitivity, can be large) [45], [27]. All of them use heuristic-based approaches. On the other hand, we start from an optimization perspective, and aim to select the threshold that can approximately minimize the optimization objective.

Using Shuffle-DP to Help LDP. Another related line of research is the shuffer-DP problem. When many user reports are anonymized and then mixed (shuffled), one can argue a stronger privacy guarantee [20], [5], [12]. Such a privacy amplification effect holds only when the anonymization party is trusted. This extension can be applied in our setting. In particular, reports of the first m values can be sent to a shuffler. Assuming these values are perturbed using an  $\epsilon$ -LDP protocol before sending them to the shuffler, it is proved in [5] that the output of those shuffled reports will satisfy  $(\epsilon', \delta)$ -DP, for some  $\epsilon' = O((1 \wedge \epsilon)e^{\epsilon} \sqrt{\log(1/\delta)/m})$ .

## VI. CONCLUSION AND DISCUSSION

We have presented a privacy-preserving algorithm ToPS to continually outsource a stream of observations. ToPS first finds a threshold to truncate the stream using the exponential mechanism, optimizing errors due to noise and bias. Then ToPS runs an on-line hierarchical structure and adds noise to the stream to satisfy differential privacy (DP). Finally, the noisy data are smoothed to improve utility. We also design a mechanism ToPL that satisfies the local version of DP. Our mechanisms can be applied to real-world applications where continuous monitoring and reporting of statistics, e.g., smart meter data and taxi fares, are required. The design of ToPS and ToPL are flexible and have the potential to be extended to incorporate more properties of the data. We list some as follows.

Shorten the Holdout of the Stream. We follow the setting of PAK [33] and use the first m values to output the threshold  $\theta$ . If we want to start outputting the stream sooner, we can use our Threshold optimizer with only a few observations to find a rough threshold. During the process of outputting the stream, we can use sequential composition (described in Section II-D) with a smaller privacy budget but more data points to fine tune the threshold.

**Update**  $\theta$ **.** We follow the setting of PAK and assume the distribution stays the same. If the distribution changes, we can have the Threshold optimizer run multiple times (using either sequential composition to update  $\theta$  while simultaneously outputting the stream, or the parallel composition theorem to block some values to update  $\theta$ ).

Utilizing Patterns of the Data. If there is further information, such that the data changes slowly (e.g., the current value and the next one differ in only a small amount), or the

data changes regularly (e.g., if the values show some Diurnal patterns), are given, we can potentially utilize that to improve the performance of our method as well.

## REFERENCES

- [1] Frequent itemset mining dataset repository. http://fimi.ua.ac.be/data/.
- [2] New york taxi trip record data. https://www1.nyc.gov/site/tlc/about/ tlc-trip-record-data.page.
- [3] J. Acharya, Z. Sun, and H. Zhang. Hadamard response: Estimating distributions privately, efficiently, and with little communication. In AISTATS, 2019.
- [4] G. Acs and C. Castelluccia. A case study: Privacy preserving release of spatio-temporal density in paris. In Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 1679–1688, 2014.
- [5] B. Balle, J. Bell, A. Gascón, and K. Nissim. The privacy blanket of the shuffle model. In *Annual International Cryptology Conference*, pages 638–667. Springer, 2019.
- [6] R. Bassily, K. Nissim, U. Stemmer, and A. G. Thakurta. Practical locally private heavy hitters. In NIPS, 2017.
- [7] R. Bassily and A. D. Smith. Local, private, efficient protocols for succinct histograms. In STOC, 2015.
- [8] J. Bolot, N. Fawaz, S. Muthukrishnan, A. Nikolov, and N. Taft. Private decayed predicate sums on streams. In *Proceedings of the 16th International Conference on Database Theory*, pages 284–295. ACM, 2013.
- [9] T.-H. H. Chan, E. Shi, and D. Song. Private and continual release of statistics. ACM Transactions on Information and System Security (TISSEC), 14(3):1–24, 2011.
- [10] V. Chandola, A. Banerjee, and V. Kumar. Anomaly detection: A survey. ACM computing surveys (CSUR), 41(3):1–58, 2009.
- [11] Y. Chen, A. Machanavajjhala, M. Hay, and G. Miklau. Pegasus: Dataadaptive differentially private stream processing. In *Proceedings of* the 2017 ACM SIGSAC Conference on Computer and Communications Security, pages 1375–1388. ACM, 2017.
- [12] A. Cheu, A. D. Smith, J. Ullman, D. Zeber, and M. Zhilyaev. Distributed differential privacy via shuffling. In *EUROCRYPT*, pages 375–403, 2019.
- [13] G. Cormode, S. Jha, T. Kulkarni, N. Li, D. Srivastava, and T. Wang. Privacy at scale: Local differential privacy in practice. In SIGMOD, 2018.
- [14] G. Cormode, T. Kulkarni, and D. Srivastava. Answering range queries under local differential privacy. PVLDB, 2019.
- [15] W.-Y. Day, N. Li, and M. Lyu. Publishing graph degree distribution with node differential privacy. In *Proceedings of the 2016 International Conference on Management of Data*, pages 123–138, 2016.
- [16] J. C. Duchi, M. I. Jordan, and M. J. Wainwright. Local privacy and statistical minimax rates. In FOCS, 2013.
- [17] J. C. Duchi, M. I. Jordan, and M. J. Wainwright. Minimax optimal procedures for locally private estimation. *Journal of the American Statistical Association*, 113(521):182–201, 2018.
- [18] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In TCC, 2006.
- [19] C. Dwork, M. Naor, T. Pitassi, and G. N. Rothblum. Differential privacy under continual observation. In *Proceedings of the forty-second ACM* symposium on Theory of computing, pages 715–724, 2010.
- [20] Ú. Erlingsson, V. Feldman, I. Mironov, A. Raghunathan, K. Talwar, and A. Thakurta. Amplification by shuffling: From local to central differential privacy via anonymity. arXiv preprint arXiv:1811.12469, 2018.
- [21] Ú. Erlingsson, V. Pihur, and A. Korolova. RAPPOR: randomized aggregatable privacy-preserving ordinal response. In *CCS*, 2014.
- [22] L. Fan and L. Xiong. An adaptive approach to real-time aggregate monitoring with differential privacy. *IEEE Transactions on knowledge* and data engineering, 26(9):2094–2106, 2013.
- [23] F. Fioretto and P. Van Hentenryck. Optstream: releasing time series privately. *Journal of Artificial Intelligence Research*, 65:423–456, 2019.

- [24] M. Hay, V. Rastogi, G. Miklau, and D. Suciu. Boosting the accuracy of differentially private histograms through consistency. *PVLDB*, 3(1), 2010
- [25] M. Joseph, A. Roth, J. Ullman, and B. Waggoner. Local differential privacy for evolving data. 2018.
- [26] G. Kellaris, S. Papadopoulos, X. Xiao, and D. Papadias. Differentially private event sequences over infinite streams. *Proceedings of the VLDB Endowment*, 7(12):1155–1166, 2014.
- [27] I. Kotsogiannis, Y. Tao, X. He, M. Fanaeepour, A. Machanavajjhala, M. Hay, and G. Miklau. Privatesql: a differentially private sql query engine. *Proceedings of the VLDB Endowment*, 12(11):1371–1384, 2019.
- [28] N. Li, W. Qardaji, D. Su, and J. Cao. Privbasis: Frequent itemset mining with differential privacy. *Proceedings of the VLDB Endowment*, 5(11):1340–1351, 2012.
- [29] Z. Li, T. Wang, M. Lopuhaä-Zwakenberg, B. Skoric, and N. Li. Estimating numerical distributions under local differential privacy. In SIGMOD, 2020.
- [30] F. McSherry and K. Talwar. Mechanism design via differential privacy. In FOCS, volume 7, pages 94–103, 2007.
- [31] A. Molina-Markham, P. Shenoy, K. Fu, E. Cecchet, and D. Irwin. Private memoirs of a smart meter. In Proceedings of the 2nd ACM workshop on embedded sensing systems for energy-efficiency in building, pages 61–66, 2010.
- [32] K. Nissim, S. Raskhodnikova, and A. Smith. Smooth sensitivity and sampling in private data analysis. In *Proceedings of the thirty-ninth* annual ACM symposium on Theory of computing, pages 75–84. ACM, 2007.
- [33] V. Perrier, H. J. Asghar, and D. Kaafar. Private continual release of real-valued data streams. ndss, 2019.
- [34] W. H. Qardaji, W. Yang, and N. Li. Understanding hierarchical methods for differentially private histograms. PVLDB, 6(14), 2013.
- [35] 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 CCS, 2016
- [36] V. Rastogi and S. Nath. Differentially private aggregation of distributed time-series with transformation and encryption. In *Proceedings of the* 2010 ACM SIGMOD International Conference on Management of data, pages 735–746, 2010.
- [37] N. Wang, X. Xiao, Y. Yang, T. D. Hoang, H. Shin, J. Shin, and G. Yu. Privtrie: Effective frequent term discovery under local differential privacy. In *ICDE*, 2018.
- [38] 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 *Proceedings of IEEE ICDE*, 2019.
- [39] Q. Wang, Y. Zhang, X. Lu, Z. Wang, Z. Qin, and K. Ren. Real-time and spatio-temporal crowd-sourced social network data publishing with differential privacy. *IEEE Transactions on Dependable and Secure Computing*, 15(4):591–606, 2016.
- [40] T. Wang, J. Blocki, N. Li, and S. Jha. Locally differentially private protocols for frequency estimation. In USENIX Security, 2017.
- [41] T. Wang, B. Ding, J. Zhou, C. Hong, Z. Huang, N. Li, and S. Jha. Answering multi-dimensional analytical queries under local differential privacy. In SIGMOD, 2019.
- [42] T. Wang, N. Li, and S. Jha. Locally differentially private frequent itemset mining. In SP, 2018.
- [43] T. Wang, N. Li, and S. Jha. Locally differentially private heavy hitter identification. *IEEE Transactions on Dependable and Secure* Computing, 2019.
- [44] S. L. Warner. Randomized response: A survey technique for eliminating evasive answer bias. *Journal of the American Statistical Association*, 60(309), 1965.
- [45] R. J. Wilson, C. Y. Zhang, W. Lam, D. Desfontaines, D. Simmons-Marengo, and B. Gipson. Differentially private sql with bounded user contribution. In *International Symposium on Privacy Enhancing Technologies Symposium*. Springer, 2020.
- [46] M. Ye and A. Barg. Optimal schemes for discrete distribution estimation under locally differential privacy. *IEEE Transactions on Information Theory*, 2018.

- [47] C. Zeng, J. F. Naughton, and J.-Y. Cai. On differentially private frequent itemset mining. *Proceedings of the VLDB Endowment*, 6(1):25–36, 2012
- [48] Z. Zhang, T. Wang, N. Li, S. He, and J. Chen. Calm: Consistent adaptive local marginal for marginal release under local differential privacy. In CCS, 2018.
- [49] Z. Zheng, R. Kohavi, and L. Mason. Real world performance of association rule algorithms. In Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining, pages 401–406. ACM, 2001.

#### APPENDIX

## A. Mechanisms of Differential Privacy

We review the method of smooth sensitivity that PAK [33] uses for estimating the percentile.

**Smooth Sensitivity.** Rather than the global sensitivity that considers any pair of neighboring sequences, the local sensitivity fixes one dataset V and only considers all possible neighboring sequence V' around V.

$$LS_{V}(f) = \max_{V': V' \cong V} ||f(V) - f(V')||_{1}$$

The advantage of using local sensitivity is that we only need to consider neighbors of V which could result in lower sensitivity of the function f, and consequently lower noise added to the true answer f. Unfortunately, replacing the global sensitivity with local sensitivity directly (e.g., in the Laplace mechanism) violates DP. This is handled by using smooth sensitivity [32] instead.

For b > 0, the b-smooth sensitivity of f at  $V \in \mathcal{V}$ , denoted by  $\mathsf{SS}_{V,b}(f)$ , is defined as

$$\mathsf{SS}_{V,b}(f) = \max_{V'} \left\{ \mathsf{LS}_{V'}(f) \cdot e^{-b \cdot d(V,V')} \right\}$$

where d(V, V') denotes the hamming distance between V and V'. The method of smooth sensitivity is given below:

$$\mathbf{A}(V) = f(V) + \frac{\mathsf{SS}_{V,b}}{a} \cdot Z$$

where Z is a random variable from some specified distribution. To obtain  $(\epsilon,0)$ -DP, Z is drawn from the Cauchy distribution with density  $\propto \frac{1}{1+|z|^{\gamma}}$  for  $\gamma>1$ . But to use an exponentially decaying distribution (which gives good accuracy), such as standard Laplace or Gaussian distribution, one can only obtain  $(\epsilon,\delta)$ -DP. For both, we set  $b\leq \frac{\epsilon}{-2\log(\delta)}$  as the smoothing parameter. If we use the Laplace distribution with scale  $1,a=\frac{\epsilon}{2}$ . When the noise is standard Gaussian, then  $a=\frac{\epsilon}{\sqrt{-\ln\delta}}$  [32].

## B. More Details about PAK

PAK computes the smooth sensitivity of the empirical p-quantile, i.e.,  $\hat{x}_p$ , as

$$\begin{split} \mathsf{SS}_{V,b}(\hat{x}_p) &= \max_{k=0,1,...,m+1} \\ \left\{ e^{-bk} \cdot \max_{t=0,1,...,k+1} [V^s(P+t) - V^s(P+t-k-1)] \right\} \end{split}$$

Here,  $V^s$  is the sorted string of the first m values of V in ascending order, where  $V^s(i) = 0$  if i < 1 and B if i > 1

m. And P is the rank of  $\hat{x}_p$ . After computing the smooth sensitivity, [32] sets the threshold  $\theta$  as

$$\theta = \hat{x}_p + \frac{\mathsf{SS}_{V,b}(\hat{x}_p)}{a} \cdot Z$$

where Z is a random variable from some specified distribution in order to satisfy DP.

In [33], the authors propose to bound  $\Pr[\theta < x_p]$  to be arbitrarily small (by an arbitrary  $\beta$ ) and thus uses

$$\theta = \hat{x}_p + \frac{\kappa SS_{V,b}(\hat{x}_p)}{a} \cdot (Z + G_{ns}^{-1}(1 - \beta)),$$

where  $\kappa$  is a positive real number  $\left(1-\frac{(e^b-1)G_{\rm ns}^{-1}(1-\beta_{\rm lt})}{a}\right)^{-1}$  and  $G_{\rm ns}$  denotes the CDF of the distribution of Z.

The threshold  $\theta$  released via the above mechanism is differentially private since  $\kappa SS_{V,b}(\hat{x}_p)$  is a smooth upper bound of  $\hat{x}_p$  and  $\kappa$  only depends on public parameters.

In their evaluation, the authors aim to get the 99.5-percentile and set p = 99.575,  $\beta = 0.3 \cdot 0.02$ , and  $\delta = 1/n^2$ .

## C. Mechanisms of Local Differential Privacy

In this subsection, we review the primitives proposed for LDP. We use v to denote the user's private value, and y as the user's report that satisfies LDP. In this section, following the notations in the LDP literature, we use p and q to denote probabilities.

1) LDP Mechanism for Density Estimation: One basic primitive of LDP is to estimate the histogram (or density if the domain is continuous). It is like what Laplace mechanism solves in the centralized DP setting.

**Square Wave.** In [29], an LDP method that can give the full density estimation is proposed. The intuition behind this approach is to try to increase the probability that a noisy reported value carries meaningful information about the input. Intuitively, if the reported value is the true value, then the report is a "useful signal", as it conveys the extract correct information about the true input. If the reported value is not the true value, the report is in some sense noise that needs to be removed. Exploiting the ordinal nature of the domain, a report that is different from but close to the true value v also carries useful information about the distribution. Therefore, given input v, we can report values closer to v with a higher probability than values that are farther away from v. The reporting probability looks like a squared wave, so the authors call the method Square Wave method (SW for short).

Without loss of generality, we assume values are in the domain of [0,1]. To handle an arbitrary range  $[\ell,r]$ , each user first transforms v into  $\frac{v-\ell}{r-\ell}$  (mapping  $[\ell,r]$  to [0,1]); and the estimated result is transformed back. Define the "closeness" measure  $b=\frac{\epsilon e^{\epsilon}-e^{\epsilon}+1}{2e^{\epsilon}(e^{\epsilon}-1-\epsilon)}$ , the Square Wave mechanism SW is defined as:

$$\forall y \in [-b,1+b], \; \Pr\left[\mathrm{SW}(v) = y\right] = \left\{ \begin{array}{ll} p, & \text{if } |v-y| \leq b \\ q, & \text{otherwise} \end{array} \right..$$

By maximizing the difference between p and q while satisfying the total probability adds up to 1, we can derive  $p=\frac{e^\epsilon}{2be^\epsilon+1}$  and  $q=\frac{1}{2be^\epsilon+1}$ .

After receiving perturbed reports from all users, the server runs the Expectation Maximization algorithm to find an estimated density distribution that maximizes the expectation of observing the output. Additionally, the server applies a special smoothing requirement to the Expectation Maximization algorithm to avoid overfitting.

2) Candidate Methods for the Perturber: As we are essentially interested in estimating the sum over time, the following methods that estimate mean within a population are useful. We first describe two basic methods. Then we describe a method that adaptively uses these two to get a better accuracy in all cases. Our perturber will use the final method.

Stochastic Rounding. This method uses stochastic rounding to estimates the mean of a continuous/ordinal domain [17]. We call it Stochastic Rounding (SR for short). Assume the private input value v is in the range of [-1,1] (otherwise, we can first projected the domain into [-1,1]), the main idea is to round v to v' so that v'=1 with probability  $p_1=\frac{1}{2}+\frac{v}{2}$  and v'=-1 w/p  $1-p_1$ . This stochastic rounding step is unbiased in that  $\mathbb{E}\left[v'\right]=v$ . Then given a value  $v'\in\{-1,1\}$ , the method runs binary random response to perturb v' into y. In particular, let  $p=\frac{e^\epsilon}{e^\epsilon+1}$  and  $q=1-p=\frac{1}{e^\epsilon+1}, y=v'$  w/p p, and p=10. The method has variance  $\left(\frac{e^\epsilon+1}{e^\epsilon-1}\right)^2-v^2$ . Piecewise Mechanism. This method is proposed in [38]. It is also used for mean estimation, but can get more accurate mean estimation than SR when e0. 1.29. In this method, the input domain is e1.31, and the output domain is e2.41, where e3.42, where e4.41, and the output domain is e5.43, where e6.42, and e6.44, and e6.45, where e7.45, and the output domain is e7.47, where

 $[\ell(v),r(v)]$  where  $\ell(v)=rac{e^{\epsilon/2}\cdot v-1}{e^{\epsilon/2}-1}$  and  $r(v)=rac{e^{\epsilon/2}\cdot v+1}{e^{\epsilon/2}-1}$ , such that with input v, a value in the range  $[\ell(v),r(v)]$  will be reported with higher probability than a value outside the range. The high-probability range looks like a "piece" above the true value, so the authors call the method Piecewise Mechanism (PM for short). The perturbation function is defined as

$$\forall_{y \in [-s,s]} \Pr\left[ \mathrm{PM}(v) = y \right] = \left\{ \begin{array}{ll} p = \frac{e^{\epsilon/2}}{2}z, & \text{if } y \in [\ell(v), r(v)] \\ q = \frac{1}{2e^{\epsilon/2}}z, & \text{otherwise} \end{array} \right.$$

where  $z=\frac{e^{\epsilon/2}-1}{e^{\epsilon/2}+1}$ . Compared to SR, this method has a variance of  $\frac{v^2}{e^{\epsilon/2}-1}+\frac{e^{\epsilon/2}+3}{3(e^{\epsilon/2}-1)^2}$  [38].

**Hybrid Mechanism.** Both SR and PM incurs a variance that depend on the true value, but in the opposite direction. In particular, when  $v=\pm 1$ , the variance of SR is lowest, but the variance of PM is highest. The authors of [38] thus propose a method called Hybrid Mechanism (HM for short) to achieve a good accuracy for any v. In particular, define  $\alpha=1-e^{-\epsilon/2}$ , when  $\epsilon>0.61$ , users use PM w/p  $\alpha$  and SR w/p  $1-\alpha$ . When  $\epsilon\leq0.61$ , only SR will be called. It is proved in [38] HM gives better accuracy than SR and PM. In particular, the worst case variance is

$$\operatorname{Var}\left[\tilde{v}\right] = \begin{cases} \left(\frac{e^{\epsilon} + 1}{e^{\epsilon} - 1}\right)^{2}, & \text{when } \epsilon \leq 0.61\\ \frac{1}{e^{\epsilon/2}} \left[\left(\frac{e^{\epsilon} + 1}{e^{\epsilon} - 1}\right)^{2} + \frac{e^{\epsilon/2} + 3}{3(e^{\epsilon/2} - 1)}\right], & \text{when } \epsilon > 0.61 \end{cases}$$
(13)