
LINKEDIN'S AUDIENCE ENGAGEMENTS API: A PRIVACY PRESERVING DATA ANALYTICS SYSTEM AT SCALE

RYAN ROGERS, SUBBU SUBRAMANIAM, SEAN PENG, DAVID DURFEE, SEUNGHYUN LEE,
SANTOSH KUMAR KANCHA, SHRADDHA SAHAY, AND PARVEZ AHAMMAD

All authors: LinkedIn

e-mail address: {rrogers,ssubramaniam,speng,ddurfee,snlee,sakumarkancha,ssahay,pahammad}@linkedin.com

ABSTRACT. We present a privacy system that leverages differential privacy to protect LinkedIn members' data while also providing audience engagement insights to enable marketing analytics related applications. We detail the differentially private algorithms and other privacy safeguards used to provide results that can be used with existing real-time data analytics platforms, specifically with the open sourced Pinot system. Our privacy system provides user-level privacy guarantees. As part of our privacy system, we include a budget management service that enforces a strict differential privacy budget on the returned results to the analyst. This budget management service brings together the latest research in differential privacy into a product to maintain utility given a fixed differential privacy budget.

1. INTRODUCTION

LinkedIn's Audience Engagement API is a platform that enables marketers (analysts) aggregated insights about members' content engagements while ensuring member (user) data is protected. Consider an advertiser that is selling a cloud solution and wants to create a sponsored post on LinkedIn. The advertiser might use the Audience Engagement API to do research and find that the target audience engages with GDPR articles. Hence, the advertiser should write about how their cloud solution adheres to GDPR standards, thus increasing engagement. By design, the Audience Engagement API is secure, aggregated, and uses state of the art differentially private algorithms to provide rigorous privacy guarantees. The data honors user privacy settings and applicable privacy laws. Further, data is purged within 30 days of a user leaving the ecosystem because it has a 30 day retention.

The primary reason to leverage differential privacy is due to *differencing attacks*, where the difference between two queries reveals an individual's content. For example, one query can be for the top articles by engagement from CEOs in India and then another query asks

Key words and phrases: differential privacy, API, data analytics, scale.

for the top articles by engagement from CEOs in India or LinkedIn. This attack makes aggregation and thresholding approaches insufficient — two results having counts above a threshold does not mean that their difference cannot uniquely identify an individual. Rather than limiting the scope of the Audience Engagement API by reducing the ways the dataset can be sliced in an ad hoc way, we instead worked to include differentially private algorithms to prevent such differencing attacks. To make such a product private, it is apparent that one must add noise, to prevent differencing attacks, and limit the number of accesses to the API, to prevent reconstructing the dataset, despite the noise that is added. Differential privacy then formalizes these approaches via randomized algorithms and its composition properties.

To incorporate differential privacy, we carefully balance various resources, including data storage distributed across several servers, real-time query computation, privacy loss quantified by the differential privacy parameters (ϵ, δ) , and accuracy. We describe here the overall privacy system deployed at LinkedIn that balances these resources to provide a product that surfaces audience engagement insights while putting members first by safeguarding their data.

Providing scalable, real-time analytics with low latency without differential privacy is challenging enough. Luckily, we have the open source real-time distributed OLAP datastore, called Apache Pinot (incubating) [1]. Pinot enables use cases like *Job and Publisher Analytics* and *Who Viewed My Profile*. In order to develop a differentially private system, we need to think how it can be used in conjunction with a (distributed) OLAP system such as Pinot. This would enable us to have scalable privacy systems. Furthermore, we need to implement a budgeting tool into the API so that analysts cannot repeatedly query the dataset thus making noise addition pointless. Our goal is twofold: implement differentially private algorithms that can be used with real-time distributed OLAP systems and incorporate a privacy budget management service to restrict the amount of information an analyst can retrieve. For the privacy budget management service, we incorporate the latest composition bounds for our particular algorithms [2] to extract more utility subject to a given differential privacy budget.

1.1. Contributions. We make several contributions toward making practical privacy systems that leverage differential privacy.

- We describe a suite of differentially private algorithms that cover the data analytics tasks for LinkedIn’s Audience Engagement API, which provide user-level privacy guarantees.
- We detail our privacy budget management service that is able to track each analyst’s privacy budget over multiple data centers. Hence, we can ensure the budget is enforced across large scale systems in real-time.
- We showcase empirical results of our algorithms on LinkedIn’s data for various privacy parameters on our deployed system.
- We provide a discussion about the considerations in our privacy system, in particular how we rationalize certain parameters. We hope that this discussion will help guide practitioners in how parameters might be set and provide transparency into our deployed system.

Although the private algorithms and privacy budget formulas were known in prior work, the main contribution of this work is in combining both algorithms and budget management into a system that can easily scale to large datasets and multiple analysts querying the system while applying the privacy system in the Audience Engagement API product at

LinkedIn. In particular, we developed a library of private algorithms and a privacy budget management system separately so that each could scale according to their own requirements; see Section 4 for more detail. Further, we propose two units of budget, *information* and *call* budgets, that can be deducted for each analyst depending on each result she receives. We can then use the latest, state of the art privacy composition formulas that tightly bound the overall privacy loss. We also state our assumptions for the privacy system in Section 8, including analysts not colluding and data churn for refreshing privacy budgets.

1.2. Related Work. Differential privacy has become the standard privacy benchmark for data analytics on sensitive datasets. Despite its popularity in the academic literature, the number of actually implemented differential privacy systems is limited, but growing. Several of the currently implemented systems with differential privacy are in the local model, where data is individually privatized prior to being aggregated on a central server. The main local differentially private systems include Google’s RAPPOR on their Chrome browser [3], Apple’s iOS and MacOS diagnostics [4], and Microsoft’s telemetry data in Windows 10 Fall Creators Update [5].

The privacy model we are interested in for this work is the global privacy setting, where data is already stored centrally, but we want to ensure each result computed on the data is privatized. In this less restrictive privacy setting, the main industrial differential privacy systems include Microsoft’s PINQ [6], Uber’s FLEX for its internal analytics [7], LinkedIn’s PriPeARL for its ad analytics [8], Google’s recent differential privacy open source project [9, 10], and the 2020 U.S. Census [11]. In this work, we present a privacy system that incorporates a privacy budget management service to ensure user-level privacy, whereas LinkedIn’s PriPeARL system provides event-level privacy and was focused on providing consistent results, which we also incorporate. The FLEX system points out that a privacy budget management service can be implemented but does not provide a strategy for how to do it. Further, our system is part of an API that allows for adaptively chosen queries computed in real-time, which is, to our knowledge, a different model from the future U.S. Census Bureau’s system.

The main difference between the approach recently proposed in [10] and this work is that we do not bound user contributions across and within different partitions.¹ Such an approach would create a significant bottleneck in processing queries in a real-time system, since each online query can require a pre-processing step over the dataset to bound user contributions. For Audience Engagement, we are dealing with terabytes of data. The **UnkGumb** algorithm provides user-level privacy guarantees for count distinct queries without pre-processing, thus handling similar *queries-per-second* (QPS) as without privacy. See Section 3 for more detail. Although Wilson et al. [10] do discuss a privacy budget, it does not consider optimized privacy loss bounds for the data analytics tasks we consider. Our system takes into account the various privacy algorithms to take advantage of the state of the art privacy composition bounds, such as *pay-what-you-get* composition and improved composition bounds for exponential mechanisms [12, 2].

There are other open source libraries for differentially private algorithms, such as PrivateSQL [13] and the recent collaboration project between Harvard’s IQSS and Microsoft [14]. The former work generates a synthetic dataset, *private synopses*, that is based on all

¹Note that a caveat in the Google open-source code is that the “implementation assumes that each user contributes only a single row to each partition.”, <https://github.com/google/differential-privacy>

queries that are posed in advance. Such an approach is very appealing, but would not be feasible in our setting due to the size of the underlying dataset and the set of all possible queries that can be asked by an analyst also being large.

Another related privacy system is PSI (Ψ) from the Harvard Privacy Tools Project [15]. PSI is a private data sharing interface to “enable researchers in the social sciences and other fields to share and explore privacy-sensitive datasets with the strong privacy protections of differential privacy.” Although they support several commonly used statistics, our system covers the necessary algorithms to privatize queries in the Audience Engagement API. Further, our system allows for handling highly distributed datasets via Pinot while enforcing a strict privacy budget that is eventually consistent across data centers.

2. PRELIMINARIES

We now present some notation and fundamental definitions that will be used to describe our privacy system. We will denote the data histogram as $\mathbf{h} \in \mathbb{N}^d$ where d is the dimension of the data universe, which might be unknown or known. We say that \mathbf{h} and \mathbf{h}' are neighbors, sometimes denoted as $\mathbf{h} \sim \mathbf{h}'$, if they differ in the presence or absence of at most one member’s data. We now define differential privacy [16, 17].

Definition 2.1 Differential Privacy. *An algorithm \mathcal{M} that takes a histogram in \mathbb{N}^d to some arbitrary outcome set \mathcal{Y} is (ϵ, δ) -differentially private (DP) if for all neighbors \mathbf{h}, \mathbf{h}' and for all outcome sets $S \subseteq \mathcal{Y}$, we have $\Pr[\mathcal{M}(\mathbf{h}) \in S] \leq e^\epsilon \Pr[\mathcal{M}(\mathbf{h}') \in S] + \delta$. If $\delta = 0$, then we simply write ϵ -DP.*

In our algorithms, we will add noise to the histogram counts. The noise distributions we consider are from a Gumbel distribution where $\text{Gumbel}(b)$ has PDF $p_{\text{Gumbel}}(z; b)$ or a Laplace distribution where $\text{Lap}(b)$ has PDF $p_{\text{Lap}}(z; b)$, and

$$p_{\text{Gumbel}}(z; b) = \frac{1}{b} \cdot e^{-(z/b + e^{-z/b})}$$

$$p_{\text{Lap}}(z; b) = \frac{1}{2b} \cdot e^{-|z|/b}.$$

As an analyst interacts with private algorithms, the resulting privacy parameters increase with each returned result. Hence, we need to account for the overall *privacy budget* that an analyst can exhaust before the privacy loss is deemed to be too large. We then use the composition property of DP to bound the resulting privacy parameters. We will use *bounded range* in our composition analysis, which was introduced by Durfee and Rogers [12]. Note that ϵ -BR mechanisms are ϵ -DP and ϵ -DP mechanisms are 2ϵ -BR.

Definition 2.2 Bounded Range. *Given a mechanism \mathcal{M} that takes a histogram in \mathbb{N}^d to outcome set \mathcal{Y} , we say that \mathcal{M} is ϵ -bounded range (BR) if for any $y_1, y_2 \in \mathcal{Y}$ and any neighboring databases \mathbf{h}, \mathbf{h}' we have*

$$\frac{\Pr[\mathcal{M}(\mathbf{h}) = y_1]}{\Pr[\mathcal{M}(\mathbf{h}') = y_1]} \leq e^\epsilon \frac{\Pr[\mathcal{M}(\mathbf{h}) = y_2]}{\Pr[\mathcal{M}(\mathbf{h}') = y_2]}$$

where we use the density function instead for continuous outcomes.

We now state the result from [2] and [18] that tightens the composition bound from Durfee and Rogers [12] which itself improved on the more general optimal DP composition bounds [19, 20].

Lemma 2.1 . *Let $\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_t$ each be ε -BR where the choice of mechanism \mathcal{M}_i at round i may depend on the previous outcomes of $\mathcal{M}_1, \dots, \mathcal{M}_{i-1}$, then the resulting composed algorithm is $(\varepsilon'(\delta), \delta)$ -DP for any $\delta \geq 0$ where $\varepsilon'(\delta)$ is the minimum of $t\varepsilon$ and*

$$\frac{t\varepsilon^2}{8} + \varepsilon\sqrt{\frac{t}{2}\ln(1/\delta)}. \quad (2.1)$$

We also can use the more complicated composition bound for BR mechanisms [2]. However we cannot use the optimal composition bound from [2] because it only applies to the non-adaptive setting. Here we are interested in the API setting which allows the user to ask adaptive queries, meaning the queries can depend on previous results.

3. PRIVATE DATA ANALYTICS

To incorporate differential privacy, we needed to consider the various tasks we want the application to handle. We will be focusing on data analytics based on histograms or counts over different domain elements. We will discuss each query type our privacy system handles, but first we need to set up some notation.

In order to provide a *user-level* privacy guarantee where all data records of a user are protected, as opposed to *event-level* where only an individual data record is protected, we consider two types of queries. The first consists of *distinct count* queries where a member can contribute a count of at most 1 to any number of elements, i.e. $\|\mathbf{h} - \mathbf{h}'\|_\infty \leq 1$ for any neighbors \mathbf{h}, \mathbf{h}' (ℓ_∞ -sensitivity). An example of such a query would be “what are the top- k articles that are shared among distinct members with a certain skill set?” The second type is *non-distinct count* queries where a member can increase the count of any element by at most $\tau \geq 1$, i.e. $\|\mathbf{h} - \mathbf{h}'\|_\infty \leq \tau$ for any neighbors \mathbf{h}, \mathbf{h}' . Note that $\tau = 1$ gives us the distinct count setting and τ can be a parameter for each non-distinct count query.

In either case, distinct count or non-distinct count queries, a member can either affect the count of an arbitrary number of elements $\|\mathbf{h} - \mathbf{h}'\|_0 \leq d$ for any neighbors $\mathbf{h}, \mathbf{h}' \in \mathbb{N}^d$ or a bounded number $\|\mathbf{h} - \mathbf{h}'\|_0 \leq \Delta$ for $\Delta < d$ (ℓ_0 -sensitivity). We separate these two cases as the *unrestricted* sensitivity and Δ -restricted sensitivity settings, respectively. In the case of unrestricted sensitivity, we will return only a fixed number of counts, say the top- k , in order to bound the privacy loss.

Scaling our privacy system across several analysts with queries that require data from multiple servers requires algorithms that can run efficiently with runtime that does not scale with the entire data domain size d . For example, for the top-10 articles engaged with by staff software engineers, we do not want to query over all articles, since there could potentially be billions of articles and would be computationally expensive and slow. For this reason, we distinguish the case where the data domain is reasonably sized and known, i.e. *known domain*, from when the data domain is very large or unknown, i.e. *unknown domain*.

We then summarize in Table 1 the set of queries that we want our privacy system to handle into *unrestricted sensitivity* or Δ -*restricted sensitivity* as well as *known domain* or *unknown domain* with the corresponding algorithms we will use for each setting. Recall that we can interpolate between distinct count queries and non-distinct count queries with the $\tau \geq 1$ parameter, so we include τ as a parameter to each of our algorithms. Furthermore, each algorithm takes a privacy parameter ε_{per} .

Restricting the ℓ_∞ -sensitivity is not part of the algorithm, rather it is done by using distinct count queries ($\tau = 1$), knowing a bound a priori, or done via a preprocessing step

	Δ -restricted sensitivity	unrestricted sensitivity
Known Domain	$\text{KnownLap}^{\Delta, \tau}$ [16]	$\text{KnownGumb}^{k, \tau}$ [21]
Unknown Domain	$\text{UnkLap}^{\Delta, \bar{d}, \tau}$	$\text{UnkGumb}^{k, \bar{d}, \tau}$

TABLE 1. DP algorithms for various data analytics tasks

on the data. For the unknown domain setting, we require a parameter \bar{d} which tells us how many elements from the original dataset that we can access in our algorithms. We think of \bar{d} as the maximum number of elements our OLAP system can return without causing significant latency. For the unrestricted sensitivity setting, we require our algorithms to return at most k elements, such as the top- k . This is due to the fact that a user can change the counts of an arbitrary number of elements. In such cases, to have any hope to bound the privacy loss, we bound the number of elements that can be returned. In Table 1, we refer to the following mechanisms: the standard Laplace mechanism from [16] is denoted as $\text{KnownLap}^{\Delta, \tau}$, which adds Laplace noise with scale proportional to $\tau/\epsilon_{\text{per}}$ to each count; the k -peeling exponential mechanism [21], which adds Gumbel noise with scale proportional to $\tau/\epsilon_{\text{per}}$ to each count and returns the elements with the top- k noisy counts, which we denote as $\text{KnownGumb}^{k, \tau}$; the generalized restricted sensitivity algorithm from [12] denoted as $\text{UnkLap}^{\Delta, \bar{d}, \tau}$, which is presented in Algorithm 3; the generalized unrestricted sensitivity algorithm from [12] denoted as $\text{UnkGumb}^{k, \bar{d}, \tau}$, which is presented in Algorithm 4.

The primary difference between querying the OLAP datastore for results with privacy as opposed to without privacy is that when querying for the top- k in the unknown domain setting, we instead fetch the top- \bar{d} and then use $\text{UnkGumb}^{k, \bar{d}, \tau}$ or $\text{UnkLap}^{\Delta, \bar{d}, \tau}$. Ideally, we would want to set $\bar{d} = d$ to get the full dataset, but that is not practical when the number of elements is large and the existing architecture potentially trims the results for efficiency. The choice of algorithm for each query can be a simple look up of the *group by* clause where if no additional information is given, then we default to the unknown domain and unrestricted sensitivity.

4. PRIVACY SYSTEM ARCHITECTURE

Existing OLAP datastores are designed to provide real-time data analytics over distributed datasets, with differential privacy not necessarily being incorporated from the beginning. Pinot is the analytics platform of choice at LinkedIn for site-facing use cases. In this section, we detail how we incorporated differential privacy with Pinot and the application. Figure 1 presents the overall system.

The application entity, based on the request received from the analyst, generates queries to the underlying database. The queries typically ask for a histogram grouped by some column. In order to apply the right algorithm for the query, the application needs to know the sensitivity and domain setting of the column as shown in Table 1. Also, the query is to be modified to fetch a potentially larger number of rows from the database.

We designed generic interfaces that are implemented by a suite of algorithms. The interfaces allow the application to:

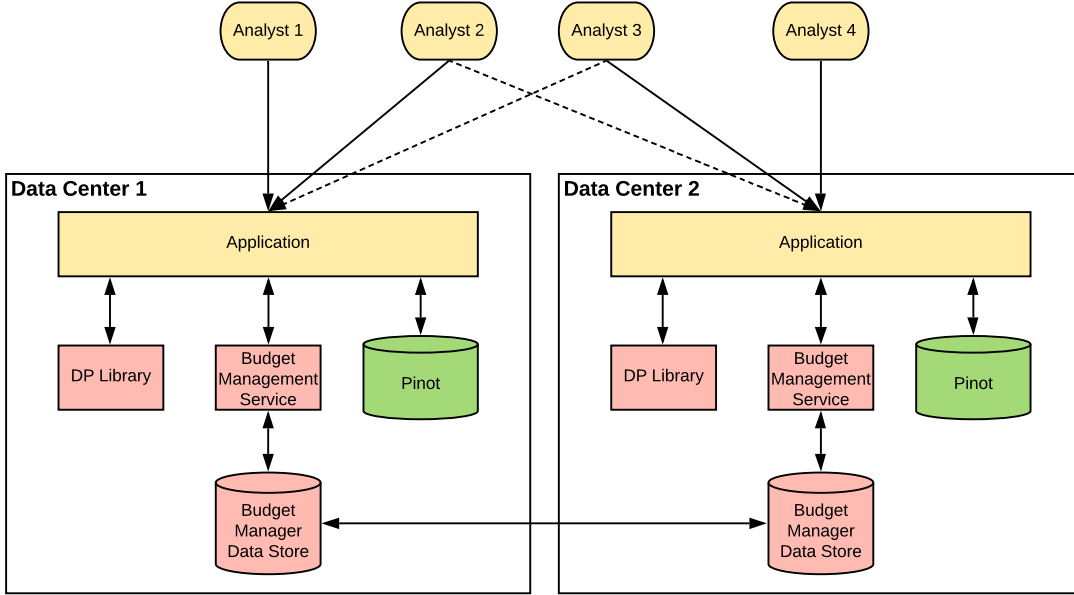


FIGURE 1. The overall privacy system with additional components for DP being the *DP Library* as well as the *Budget Management Service* and *Data Store*. The arrows between *Analysts* and *Data Centers* show that an analyst may be initially assigned one data center (bold) but can migrate to a different one (dashed).

- Retrieve modified query parameters (e.g. change k to \bar{d}).
- Estimate privacy cost of the query (e.g. return Δ or k).
- Add noise to results based on configured parameters.
- Compute the actual cost of the query, e.g. the number of items returned in the unrestricted sensitivity setting.

The application can independently invoke budget management functions, such as the following:

- Getting the available budget for an analyst to verify whether a query can even start to execute.
- Depleting the available budget with the executed query.

The cost of a query could be multi-dimensional, including the cost of making the call and of information retrieved (see Section 6).

Given the query from the analyst and the selected DP algorithm, the application will then interact with the DP library. It will first determine the expected cost of the resulting query to show the application, which is a function of the query that is asked and the selected algorithm. The application then calls the DP library to translate the query to a DP version that will be used to query Pinot. For example, if the query is for top- k and the algorithm is in the unknown domain setting, then the translation could simply modify k to $2k$, in which case $\bar{d} = 2k$ in $\text{UnkLap}^{\Delta, \bar{d}, \tau}$ and $\text{UnkGumb}^{k, \bar{d}, \tau}$. On the other hand, if the query is over the

known domain setting, then we will want to translate k to d in order to get counts over the full domain, including elements with zero counts.

Now that the application has the modified query from the DP library, we need to check whether there is enough budget remaining from the budget management service for the query to be evaluated, which are updated parameters (k^*, ℓ^*) that decrease from some fixed values. We typically call the k^* parameter the *information budget* and can be thought of as the amount of the ε parameter in DP we are consuming. Additionally, we refer to ℓ^* as the *call budget* and is associated with the δ parameter in DP. We assume that each analyst will have its own budget and each analyst starts with the same budget. If the budget is exhausted for an analyst then the budget management service does not allow the query to be executed and tells the application that the analyst has exhausted their entire budget. If the budget is not depleted, yet what remains is less than the expected cost of the query, then we still do not evaluate the query.

Once the budget management service allows for the query to be evaluated, the application queries Pinot as it would have without the privacy system only now with the translated query. The Pinot result is then returned to the application and then the application makes another call to the DP library with the Pinot result. The DP library will then run the corresponding DP algorithm on the Pinot result and return the DP result. Based on the DP result, the budget management service updates the parameters k^*, ℓ^* , as described in Algorithm 5, and returns the result to the application.

We built the algorithms module and the budget management module to be independent of each other for the following reasons:

- While DP algorithms are running on the application layer, budget management operations require a remote call to a distributed system because the budget management service needs to provide a consistent view to all application instances. Therefore, keeping the budget management independent of algorithms will allow us to scale them independently. The algorithms will need to scale to minimize memory and CPU usage, whereas the budget management service will need to scale in terms of handling higher query-per-second (QPS), while minimizing latency.
- We require that newer (as yet unknown) algorithms still be able to use and manage budgets.
- Multiple implementations of the budget manager are possible depending on system requirements. We need to be able to iterate on these independently and quickly.
- The algorithms need not (and do not) know about the analyst that is querying, and the budget manager does not behave differently depending on the type of query or algorithms used. As long as they both have a common notion of the units (k^*, ℓ^*) and dimensions of cost, it makes sense to keep these independent of each other.

4.1. Pinot: a distributed OLAP datastore. Pinot [1] is a distributed, real-time, columnar OLAP data store, currently incubating in Apache. At LinkedIn, we have two main categories of analytics applications: internal applications (such as dashboards, anomaly detection platform, A/B testing, etc.) and site-facing applications (such as *Who viewed my profile*, Talent Insights, etc.). Internal dashboards need to process a large volume of data (trillions of records), but can tolerate latencies in hundreds of milliseconds. They also have a relatively low query volume. The site-facing applications, on the other hand, serve

hundreds of millions of LinkedIn members, and therefore have a very high query volume with a latency budget of a few to perhaps tens of milliseconds.

Pinot has a flexible architecture and supports a wide variety of applications in the spectrum. Pinot production clusters at LinkedIn are serving tens of thousands queries per second, supporting more than 50 analytical use cases, and ingesting over millions of records per second. Other companies such as Uber, Microsoft, and Weibo are also operating production Pinot clusters.

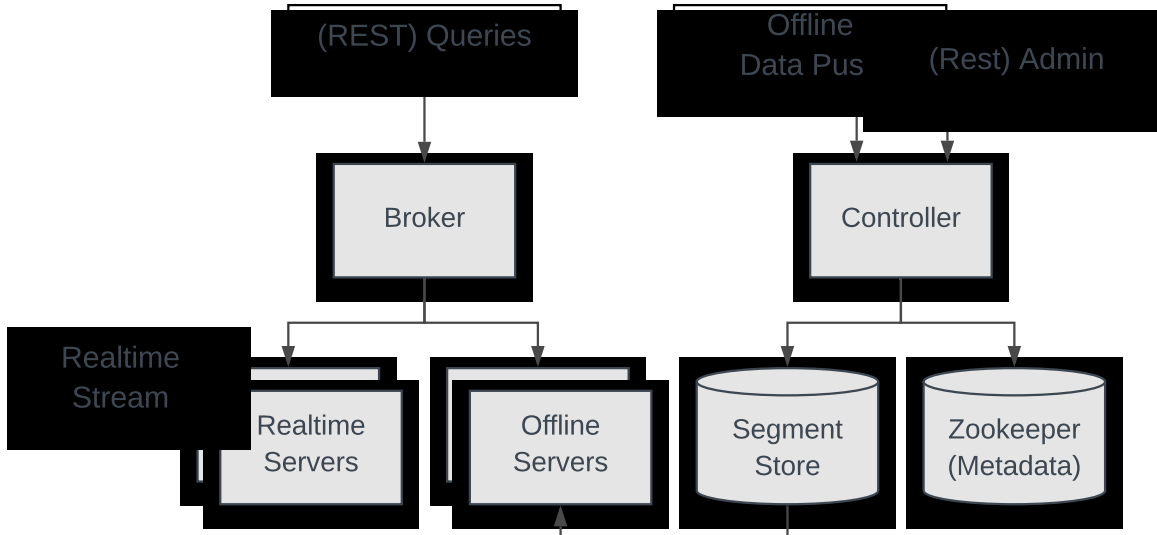


FIGURE 2. Overall Pinot Architecture.

As shown in Figure 2, Pinot has three different components: controller, broker, and server. Controllers handle cluster wide coordination, run periodic tasks for cluster state validation and retention management, and provide a REST API for managing cluster metadata. Brokers receive queries and federate them to servers so as to cover all the segments (shards) of a table. Servers execute the query on the segments. Offline servers host segments that are batch ingested while real-time servers host the segments that are ingested from streaming sources, such as Kafka [22].

For the privacy system at LinkedIn, we naturally decided to use Pinot as an OLAP data store because Pinot already supported a lot of customer-facing analytics applications like Audience Engagement. However, it is noteworthy that our architecture keeps the the budget management service and Pinot as separate components so that we can easily provide DP features to other analytical query engines such as Presto and Spark SQL.

4.2. Key-value Based Budget Management System in Espresso. We now describe our key-value based budget management system. We create one key per analyst of a table (or per use case which may have multiple tables), and the data against the key is atomically changed when we need to update the budget. The store needs to provide ways to do the read-modify-write operations, and the latency should be relatively low.

The value record will contain the following items:

- Maximum budget allowed for the user.

- The time period over which this budget is allowed (typically the time period during which the data is refreshed completely).
- The total budget used so far (or, that remain).
- The timestamp when the used budget was reset to 0 (e.g. for a monthly refresh, this will be the 1st of the month).

There are three methods that the budget manager needs to support:

- To check whether an analyst’s ID has enough budget to run a query that will consume at most a given cost, we use `checkBudget(ID, cost)`, which returns either true or false.
- To deduct an analyst’s budget, with a given ID, after getting a DP result, with a given cost, we use `updateBudget(ID, cost)`.
- To get the current budget of an analyst, with a given ID, we use `getBudget(ID)`, which returns either the analyst’s current budget that has been used or the maximum budget allowed if there has been a budget refresh.

We use an *Espresso* [23] key-value store to manage the budget. The read requests should be fast, given it is only a primary key lookup. So, a call to get the current usage (currently coming in at an unknown rate) can be fast. Espresso was chosen due to several reasons: eventual consistency in cross-datacenter replication to ensure an analyst does not exceed a given budget, capability to scale to millions of users while still keeping a fairly constant response time, control over the refresh time period, and flexibility to change per-analyst maximum upon demand.

5. DIFFERENTIALLY PRIVATE ALGORITHMS

We detail the algorithms for the various tasks in Table 1. These algorithms consist of previous work from [16], [21], and [12], or slightly modified forms. Each algorithm takes a ϵ_{per} privacy parameter, which determines the amount of noise to add, while each algorithm in the unknown domain setting has an additional $\delta > 0$ privacy parameter. We point out that `UnkGumbk, \bar{d} , τ` is the default algorithm to use when no other information is known. However, the benefit of knowing the domain is that when k results are requested, k results will be returned each time, whereas the unknown domain setting may return fewer than k . The benefit of the Δ -restricted sensitivity setting is that the budget depletes by only Δ in the known domain setting, rather than by the number of elements returned, as in the unrestricted setting.

5.1. Known Domain Algorithms. We will now state the well known Laplace [16] and Exponential [21] mechanisms. We present the Laplace mechanism [16] in the context of histogram data with the assumption that the ℓ_∞ -sensitivity between any neighbors is bounded by τ . Note that we will use a slightly different scale of noise in procedure `KnownLap Δ , τ` in Algorithm 1 than is traditionally used. This is because we want to compose bounded range algorithms in our privacy budget manager, where each algorithm has the same parameter ϵ_{per} . We go into more detail on the privacy budget service in Section 6.

We now discuss the Exponential Mechanism [21] in full generality and use the *range* of a quality score rather than the global sensitivity of the score, as was presented in [2].

Algorithm 1 $\text{KnownLap}^{\Delta, \tau}$; Laplace mechanism over known domain with ℓ_∞ -sensitivity τ , and Δ -restricted sensitivity

Input: Histogram \mathbf{h} , Δ sensitivity, along with parameter ε_{per} .

Output: Noisy histogram.

for $i \in [d]$ **do**

$$v_i = h_i + \text{Lap}(2\tau/\varepsilon_{\text{per}})$$

Return $\{v_1, \dots, v_d\}$

Definition 5.1 Exponential Mechanism. *The Exponential Mechanism $M_q : \mathcal{X} \rightarrow \mathcal{Y}$ with quality score $q : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ is written as $M_q(x)$, which samples y with probability proportional to $\exp(\varepsilon_{\text{per}}q(x, y)/S_q)$ where,*

$$S_q := \sup_{x \sim x'} \left\{ \max_{y \in \mathcal{Y}} \{q(x, y) - q(x', y)\} - \min_{y' \in \mathcal{Y}} \{q(x, y') - q(x', y')\} \right\}.$$

Note that the Exponential Mechanism is equivalent to adding Gumbel noise $\text{Gumbel}(S_q/\varepsilon_{\text{per}})$ to $q(x, y)$ for each $y \in \mathcal{Y}$ and reporting the largest noisy counts [12]. We then have the following result from [21, 2]

Lemma 5.1 . *The Exponential Mechanism is ε_{per} -BR and, hence ε_{per} -DP.*

In our case, the quality score will simply be the heights of the histogram. Note that we have only discussed the Exponential Mechanism to return a single element. In the case where we want to return k -elements, we can iteratively apply the Exponential Mechanism by removing the element that is returned in each round and then run the Exponential Mechanism again without the previously returned elements, also known as *peeling*. However, we can implement this more efficiently by adding Gumbel noise to all the counts and then releasing the top- k elements in a single shot [12]. However, we need to also include counts, so we add independent Laplace noise to the counts of the elements in the noisy top- k . We then formally present the $\text{KnownGumb}^{k, \tau}$ procedure in Algorithm 2.

Algorithm 2 $\text{KnownGumb}^{k, \tau}$; Exponential Mechanism over known domain with ℓ_∞ -sensitivity τ and unrestricted sensitivity

Input: Histogram \mathbf{h} , number of outcomes k , and parameter ε_{per} .

Output: Ordered set of indices and counts.

for $i \in [d]$ **do**

$$v_i = h_i + \text{Gumbel}(\tau/\varepsilon_{\text{per}})$$

Sort $\{v_i\}$ where $v_{i_1} \geq \dots \geq v_{i_d}$

$$\{Z_i\}_{i=1}^k \stackrel{i.i.d.}{\sim} \text{Lap}(2\tau/\varepsilon_{\text{per}})$$

Return $\{(i_1, h_{i_1} + Z_1), \dots, (i_k, h_{i_k} + Z_k)\}$

We then have the following result which follows from Dwork et al. [16], as well as from McSherry and Talwar [21, 2].

Lemma 5.2 . *Assume that $\|\mathbf{h} - \mathbf{h}'\|_\infty \leq \tau$ and $\|\mathbf{h} - \mathbf{h}'\|_0 \leq \Delta$ for any neighbors \mathbf{h}, \mathbf{h}' . The procedure $\text{KnownLap}^{\Delta, \tau}$ is $\Delta \varepsilon_{\text{per}}/2$ -DP and $\Delta \varepsilon_{\text{per}}$ -BR. Further, if Δ is large or unknown then $\text{KnownGumb}^{k, \tau}$ is $3k \varepsilon_{\text{per}}/2$ -DP and $2k \varepsilon_{\text{per}}$ -BR.*

Algorithm 3 $\text{UnkLap}^{\Delta, \bar{d}, \tau}$; Laplace mechanism over unknown domain with access to $\bar{d} + 1 > \Delta$ elements, ℓ_∞ -sensitivity τ , and Δ -restricted sensitivity.

Input: Histogram \mathbf{h} , Δ sensitivity, cut off at $\bar{d} + 1$, and $\varepsilon_{\text{per}}, \delta$.

Output: Ordered set of indices and counts.

Solve for $\hat{\delta}$: $\delta = \hat{\delta}/4 \cdot (e^{\varepsilon_{\text{per}}/2} + 1)(3 + \ln(\Delta/\hat{\delta}))$

Sort $h_{(1)} \geq h_{(2)} \geq \dots \geq h_{(\bar{d}+1)}$.

$v_\perp = h_{(\bar{d}+1)} + \tau \cdot (1 + 2\Delta \ln(\Delta/\hat{\delta})/\varepsilon_{\text{per}}) + \text{Lap}(2\tau\Delta/\varepsilon_{\text{per}})$

for $i \leq \bar{d}$ **do**

Set $v_i = h_{(i)} + \text{Lap}(2\tau\Delta/\varepsilon_{\text{per}})$

Sort $\{v_i\} \cup v_\perp$

Let v_{i_1}, \dots, v_{i_j} be the sorted list until v_\perp

Return $\{(i_1, v_{i_1}), \dots, (i_j, v_{i_j}), (\perp, v_\perp)\}$.

5.2. Unknown Domain with Δ -Restricted Sensitivity. For our unknown domain algorithms, we introduce a \perp character to denote a null element that is not part of the domain and whose count is a noisy threshold where no element with smaller noisy count is returned. We present the $\text{UnkLap}^{\Delta, \bar{d}, \tau}$ procedure in Algorithm 3 in a more general form than in [12], which only considered the distinct count case, i.e. $\tau = 1$. Further, the proof of privacy remains true if we release the counts as well as the indices. For completeness, the proof of the following result is presented in the appendix.

Lemma 5.3 Durfee and Rogers [12]. *Assume that $\|\mathbf{h} - \mathbf{h}'\|_\infty \leq \tau$ and $\|\mathbf{h} - \mathbf{h}'\|_0 \leq \Delta$ for any neighbors \mathbf{h}, \mathbf{h}' , then the procedure $\text{UnkLap}^{\Delta, \bar{d}, \tau}$ is $(\varepsilon_{\text{per}}/2, \delta)$ -DP.*

5.3. Unknown Domain with Unrestricted Sensitivity. We present the $\text{UnkGumb}^{k, \bar{d}, \tau}$ procedure in Algorithm 4 in a more general form than in [12], which only considered the distinct count case, i.e. $\tau = 1$. The proof of the following theorem follows the same analysis as in [12]. Note that we use the optimal threshold index procedure from Algorithm 6 in Durfee and Rogers [12] by default and return counts by adding Laplace noise to the discovered elements in the top- k .

Theorem 1 Durfee and Rogers [12]. *Assume $\|\mathbf{h} - \mathbf{h}'\|_\infty \leq \tau$ for any neighbors \mathbf{h}, \mathbf{h}' . Then $\text{UnkGumb}^{k, \bar{d}, \tau}$ is $((2k + 1)\varepsilon_{\text{per}}, \delta)$ -DP.*

Algorithm 4 $\text{UnkGumb}^{k, \bar{d}, \tau}$; Unknown domain mechanism with access to $\bar{d} + 1 > k$ elements, ℓ_∞ -sensitivity τ

Input: Histogram \mathbf{h} ; outcomes k , cut off at $\bar{d} + 1$, and $\varepsilon_{\text{per}}, \delta$.

Output: Ordered set of indices and counts.

Sort $h_{(1)} \geq h_{(2)} \geq \dots \geq h_{(\bar{d}+1)}$.

for $i \in \{k, \dots, \bar{d}\}$ **do**

Set $v_i = h_{(i+1)} + \tau + \tau \ln(i/\delta)/\varepsilon_{\text{per}} + \text{Gumbel}(\tau/\varepsilon_{\text{per}})$

Set $\bar{k} = \text{argmin}\{v_i\}$.

Set $h_\perp = h_{(\bar{k}+1)} + \tau \cdot (1 + \ln(\min\{\bar{k}, \bar{d} - \bar{k}\}/\delta)/\varepsilon_{\text{per}})$.

Set $v_\perp = h_\perp + \text{Gumbel}(\tau/\varepsilon_{\text{per}})$.

for $j \leq \bar{k}$ **do**

if $h_{(j)} > h_{(\bar{k}+1)}$ **then**

Set $v_{(j)} = h_{(j)} + \text{Gumbel}(\tau/\varepsilon_{\text{per}})$.

Sort $\{v_{(j)}\} \cup v_\perp$.

Let $v_{i_1}, \dots, v_{i_j}, v_\perp$ be the sorted list up until v_\perp .

$\{Z_i\}_{i=1}^j \stackrel{i.i.d.}{\sim} \text{Lap}(2\tau/\varepsilon_{\text{per}})$

if $j < k$ **then**

Return $\{(i_1, h_{i_1} + Z_1), \dots, (i_j, h_{i_j} + Z_j), \perp\}$

else

Return $\{(i_1, h_{i_1} + Z_1), \dots, (i_k, h_{i_k} + Z_k)\}$.

6. PRIVACY BUDGET MANAGEMENT SERVICE

We ultimately want to ensure that no analyst can identify any individual's data with high confidence. We then impose a strict overall $(\varepsilon^*, \delta^*)$ -DP guarantee. In order to compute the parameters $(\varepsilon_{\text{per}}, \delta)$ that we use in each call to our algorithms² over an entire sequence of interactions with the API, we also want to know how many queries the API will allow, denoted as ℓ^* that we term the *call budget*, which will effectively impact δ^* . Further, we want to track the number of elements we want to return, denoted as k^* that we term the *information budget*, which will effectively impact ε^* . Note that k^* does not necessarily equal the number of elements returned, because we might be in the restricted sensitivity setting, and ℓ^* does not precisely equal the number of calls to the API, since we might be in the known domain setting for some queries. Once we have $(\varepsilon_{\text{per}}, \delta)$, we will only use these parameters in each algorithm, hence not allowing for adaptively changing privacy parameters.

6.1. Budget Management Implementation. As mentioned in Section 4, the budget manager needs to be a distributed system so that it can be accessed/updated from different application execution platforms. Each analyst may access data from multiple data centers and each access must deduct from the same budget. Hence, the budget manager maintains eventual consistency across data centers.

The budget can be thought of as an associative array with keys from $[\ell, k]$ and values as the corresponding units used. Given a particular outcome o from the API, the budget

²Note that for the Laplace and Exponential mechanisms in the known domain, $\delta = 0$.

service will update $k^* \leftarrow k^* - \Delta$ for Δ -restricted sensitivity queries or $k^* \leftarrow k^* - 2|o|$ for unrestricted sensitivity queries where $|o|$ denotes the number of elements returned in outcome o . Furthermore, the privacy budget management system will update $\ell^* \leftarrow \ell^* - 1$ for each query the analyst makes that is in the unknown domain setting. Once k^* or ℓ^* are depleted, we prevent the analyst from making any other queries. We address the challenge of computing the individual privacy parameters $(\varepsilon_{\text{per}}, \delta)$ given $(\varepsilon^*, \delta^*, k^*, \ell^*)$ in Theorem 2.

We adopt a privacy budget management service that assumes any user does not collude with other analysts. Hence each analyst is given her own privacy budget to interact with the Audience Engagement API and her queries do not impact the budget of another analyst. One can imagine variants of this assumption, such as all analysts that belong to the same company must share a budget. Further, the API adheres to the privacy budget up to some time frame. Thus, if an analyst has asked more than ℓ^* unknown domain queries, then she will not be allowed any further queries. After this prescribed time frame, the parameters effectively get refreshed and the analyst can continue asking queries. Refresh is acceptable at regular intervals if the underlying data is flushed and replaced at similar intervals, whether through complete snapshot replacement, or rolling windows, such that the user’s data does not remain constant.

The application links with a budget manager client library so as to hide the implementation details of the budget management service because application writers do not need to know the details about the budget database, or the budget refresh mechanisms. The parameters (k^*, ℓ^*) may be configured by the application.

6.2. Differential Privacy Composition. We present pseudocode for the privacy budget management service in Algorithm 5. We then present a way to compute the privacy guarantee of our overall system, which largely follows the analysis from Durfee and Rogers [12]. Essentially, the analysis follows from the fact that each algorithm can be represented as an iterative sequence of ε_{per} -BR algorithms. Note that the algorithms in the unknown domain setting have a probability δ of larger privacy loss, which we account for in the overall δ^* in the privacy guarantee.

In order to allow for the budget management service to return counts in the unrestricted sensitivity setting, we need to account for that in our overall budget. Further, in the unknown domain/unrestricted sensitivity setting, if the last element of o_i , denoted as $o_i[-1]$, is \perp at round i then adding Laplace noise with parameter $2\tau_i/\varepsilon_{\text{per}}$ to the counts of each of the discovered $|o_i| - 1$ elements will ensure ε_{per} -BR for each count. We can then apply our privacy loss bounds to get an overall DP guarantee by updating $k^* \leftarrow k^* - 2|o_i|$ and when the last element in o_i is not \perp , then we instead update $k^* \leftarrow k^* - (2|o_i| + 1)$. Note that if we did not require counts in the results and need only return an ordered list of elements in the top- k , then we need only update $k^* \leftarrow k^* - (|o_i| + 1)$.

Theorem 2 . For $\delta' \geq 0$ and $\varepsilon_{\text{per}}, \delta > 0$, the $\text{BudgetSystem}^{k^*, \ell^*}$ is $(\varepsilon^*, \delta^*)$ -DP where $\delta^* = 2\ell^*\delta + \delta'$ and ε^* is defined as the minimum between $k^*\varepsilon_{\text{per}}$ and the following,

$$k^*\varepsilon_{\text{per}}^2/8 + \varepsilon_{\text{per}}\sqrt{\frac{k^*}{2}\ln(1/\delta')}. \quad (6.1)$$

Proof. For the Δ -restricted sensitivity setting, we are deducting the information budget by Δ in the known domain setting or we scale the privacy parameter by Δ and deduct one from the information budget in the unknown domain setting. For a given histogram in the

Algorithm 5 BudgetSystem^{k*, ℓ*}; Budget Management

Input: An adaptive stream of histograms $\mathbf{h}_1, \mathbf{h}_2, \dots$, fixed integers k^* and ℓ^* , along with per iterate privacy parameters $\varepsilon_{\text{per}}, \delta$.

Output: Sequence of outputs (o_1, o_2, \dots) .

while $k^* > 0$ and $\ell^* > 0$ **do**

Select $\mathbf{h}_i \in \mathbb{N}^{d_i}$ with ℓ_∞ -bound τ_i .

Select k_i and number of elements allowed to access \bar{d}_i

if histogram has Δ -restricted sensitivity **then**

if $\Delta > k^*$ **then**

Break

if $\bar{d}_i > d_i$ **then**

$o_i = \text{KnownLap}^{\Delta, \tau_i}(\mathbf{h}_i)$.

Update $k^* \leftarrow k^* - \Delta$.

else

$o_i = \text{UnkLap}^{\Delta, \bar{d}_i, \tau_i}(\mathbf{h}_i)$

Update $\ell^* \leftarrow \ell^* - 1$.

Update $k^* \leftarrow k^* - 1$.

if histogram has unrestricted sensitivity **then**

if $2k_i > k^*$ **then**

Break

if $\bar{d}_i > d_i$ **then**

$o_i = \text{KnownGumb}^{k_i, \tau_i}(\mathbf{h}_i)$.

Update $k^* \leftarrow k^* - 2k_i$.

else

$o_i = \text{UnkGumb}^{k_i, \bar{d}_i, \tau_i}(\mathbf{h}_i)$

Update $\ell^* \leftarrow \ell^* - 1$.

Update $k^* \leftarrow k^* - (2|o_i| + 1 - \mathbb{1}\{o_i[-1] = \perp\})$.

Return $o = (o_1, o_2, \dots)$

known domain setting, adding $\text{Lap}(2\tau/\varepsilon_{\text{per}})$ to each count will ensure $\Delta\varepsilon_{\text{per}}$ -BR. From [18], we know that each count is also $\varepsilon_{\text{per}}^2/8$ -zero-mean Concentrated DP (zCDP), defined in [24]. We can also analyze this mechanism as if we iteratively add $\text{Lap}(2\tau/\varepsilon_{\text{per}})$ to each count and then apply composition over Δ bins. We need only apply composition for the number of elements that can actually change between neighboring datasets, i.e. Δ , and not the full dimension of the histogram. For all settings, the application of each Laplace mechanism is $\varepsilon_{\text{per}}/2$ -DP and hence $\varepsilon_{\text{per}}^2/8$ -zCDP, while each application of the Exponential Mechanism is ε_{per} -BR and hence $\varepsilon_{\text{per}}^2/8$ -zCDP.

For the unknown domain algorithms, we note that each call to UnkLap is $(\varepsilon_{\text{per}}/2, \delta)$ -DP and hence δ -approximately $\varepsilon_{\text{per}}^2/8$ -zCDP. Further, the pay-what-you-get analysis of UnkGumb allows us to analyze the composition of exponential mechanisms, while also accounting for the δ terms, i.e. when events can occur in one neighboring histogram but not in another. Putting this together, the result is analyzing the composition of k^* many $\varepsilon_{\text{per}}^2/8$ -zCDP mechanisms, of which ℓ^* of them are UnkGumb or UnkLap , which is δ -approximately $\varepsilon_{\text{per}}^2/8$ -zCDP. Thus, we apply the composition bounds for zCDP from [24] and then convert to a DP claim. \square

Given the total budget for the number of outcomes and queries (k^*, ℓ^*) along with privacy budget $(\varepsilon^*, \delta^*)$ we can solve for the parameter ε_{per} that satisfies the budget, which is then used in each algorithm. One approach we can use is the following (somewhat arbitrary) choice for $\delta = \frac{\delta^*}{6\ell^*}$ and $\delta' = \delta^*/2$.

7. RESULTS

We now present some preliminary results of our privacy system for the Audience Engagement API. In Figure 3 we present curves for the number of discovered elements in a top-50 query with varying ε_{per} and \bar{d} , i.e. the number of elements to collect, in procedure $\text{UnkGumb}^{50, \bar{d}, 1}$ from Algorithm 4 with a fixed $\delta = 10^{-10}$. The query is to find the top articles that distinct members from the San Francisco area are engaging with. We provide intervals that contain the 25th and 75th percentiles over 1000 independent trials. Note that the randomness in each trial is solely from the noise generation and we are using the same dataset each time. We see that with the same level of privacy, increasing the number of elements to fetch allows us to *discover* more elements. Hence, we see a natural tradeoff not just between privacy (ε_{per}) and utility (number of elements returned), but also between run time (fetching more results) and utility. For example, we can return twice as many elements if we fetch four times more elements with Pinot and setting $\varepsilon_{\text{per}} = 0.08$.

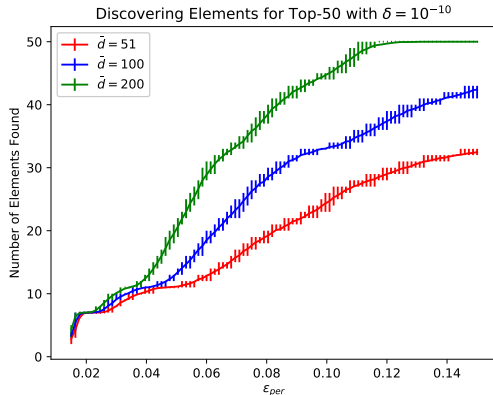


FIGURE 3. The number of returned elements in $\text{UnkGumb}^{50, \bar{d}, 1}$ for a top-50 query with various \bar{d} . We give the empirical average in 1000 trials and the (25%, 75%) percentiles.

We also empirically evaluate procedure $\text{UnkLap}^{\Delta, \bar{d}, 1}$ from Algorithm 3 in the unknown domain, Δ -restricted sensitivity setting. In Figure 4 we show both the proportion of times in 1000 trials that each element was returned (right vertical axis) as well as the comparison between the noisy counts (in green) and the true counts (in red) that are returned for the discovered elements for a single trial (left vertical axis). In each plot there is a privacy parameter $\varepsilon_{\text{per}} \in \{0.1, 0.2\}$, with fixed $\delta = 10^{-10}$. We ask for the top primary job titles of members that engaged with articles about *privacy* or *California*. We assume that any one member cannot have more than one primary job title, hence $\Delta = 1$, and fetch $\bar{d} = 1000$ results from Pinot.

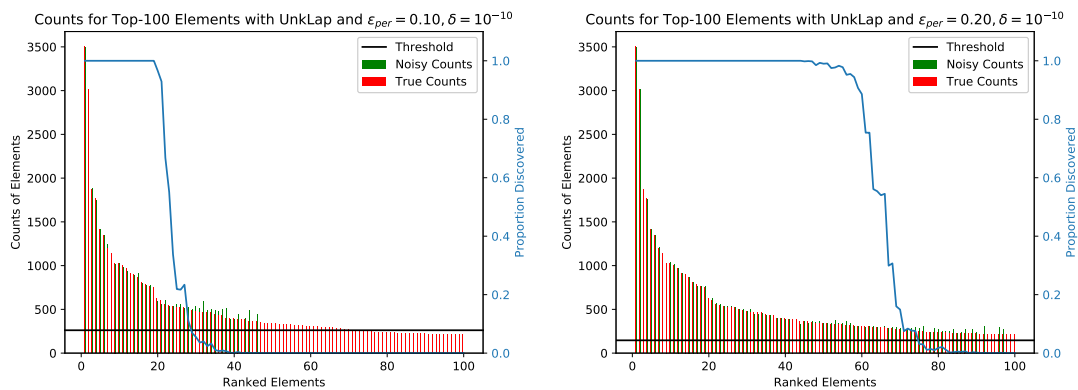


FIGURE 4. The noisy counts (left y -axis) of the discovered elements returned in $\text{UnkLap}^{1,1000,1}$ for a top-100 query as well as the proportion (right y -axis) in which various elements in 1000 independent trials were discovered. The top plot gives results for $\epsilon_{\text{per}} = 0.1$ and the bottom plot gives $\epsilon_{\text{per}} = 0.2$.

For the budget manager, we have a fixed budget for each marketing partner. Once the privacy budget is depleted, a marketing partner would recycle old queries to get the same results or wait some fixed amount of time for the privacy budget to be refreshed. This policy decision for the rate in which to refresh the budget is dependent on how often the underlying dataset gets renewed and the characteristics of the underlying dataset. In order to maintain consistency across the same queries on the same dataset, we use the same seed in the pseudorandom noise, as in [8].

8. DEPLOYMENT CONSIDERATIONS

We now discuss our approach in deploying such a system that integrated multiple components, including Pinot for data analytics, differentially private algorithms, and a privacy budget management system. Not knowing how external marketing partners would respond to budgeting access to queries and noisy results, we proceeded with a phased approach deploying our privacy system, first by turning on our privatized algorithms and only tracking usage of budget and then moving to enforce a given privacy budget. Recall that there are multiple parameters to set in our system and we detail the approach that we took to set them. Ultimately, our privacy approach was guided by developing differentially private algorithms so that privacy loss could be quantified and we focused on specific attacks for how to set parameters. Currently, the Audience Engagement API is still only available to a set of trusted users and pre-general availability (pre-GA).

8.1. Phased Approach of Deploying our Privacy System. Given the multiple teams and components that made up our privacy system, we wanted to better understand the impact to the customer when the various components were enabled. The main questions we faced included: how would external partners react with getting fewer results than they asked for (due to our private algorithms setting a data dependent threshold), and how much budget should we set without drastically modifying the behavior of how the external marketing partners interacted with the API. We then sought to answer each question separately, with

a phased approach of turning on each component. The Audience Engagement API was planned to go through a soft launch period, followed by onboarding trusted partners, to then full scale deployment. This allowed us to use the stages to deploy the different components of our privacy system.

We first deployed our differentially private algorithms and identified the various columns in the data table that had a known/unknown domain or a restricted/unrestricted sensitivity. We then tracked the privacy budget usage of the analysts that queried the API. Recall that we refer to the quantities k^* and ℓ^* as *information budget* and *call budget*, respectively. Figure 5 shows the percentage of users that would exhaust their information and call budgets over the number of days since their budget was refreshed with various cut off amounts. In this application, we refresh the budget once a month (see next subsection for more details on this), so we show the number of days since refresh on the x -axis because not all users will make a query on the first of the month and their budget does not refresh until they make their first query in the month. Further, we see that there is a percentage of users that would exhaust their budget on the first day their budget is refreshed due to asking a batch of queries at once.

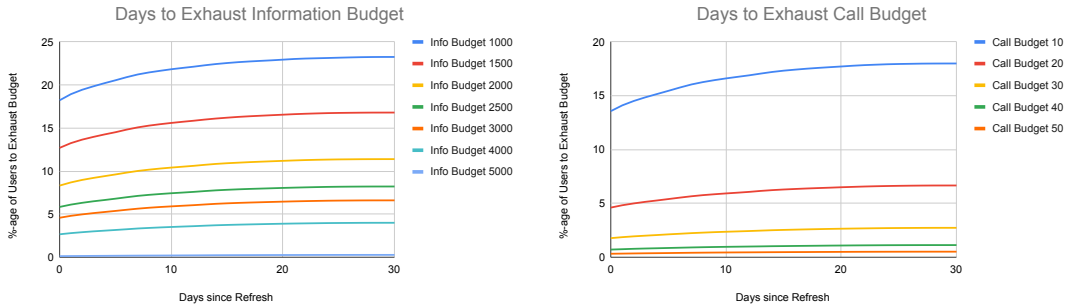


FIGURE 5. We plot the percentage of analysts that exceed their information or call budgets from the time that budget is refreshed for various information and call budgets.

The main takeaway from the plots in Figure 5 is that we can set information and call budgets in a way that does not limit a vast majority of users. In particular, setting information budget to 3000 would not impact more than 93% of users and a call budget of 30 would not impact more than 95% of users.

8.2. Consistency and Data Refresh. It is important to point out, from a product utility perspective, that data consistency is crucial. Although this might seem to contradict the inherent randomness of differential privacy, we still want to ensure that if someone asks the same query, then they get the same result. To ensure consistent results, we use the pseudorandom seed generation from [8] for our randomized algorithms. Hence, with the same query, we will use the same pseudorandom seed and hence the results will be the same, unless the underlying dataset has changed. The trusted parties who were granted access to the API all built UIs to facilitate access to their users, which prevents arbitrary query construction with small syntactic changes that leaves the semantics of the query unchanged.

This pseudorandom seed has an extra benefit for privacy as well, since if an analyst asks the same query multiple times, the random answers will not concentrate to the true answer.

Ensuring that private answers do not concentrate to the true result is one of the primary reasons for the privacy budget management, but setting information cost and call cost to ensure that the exact same query ran repeatedly does not concentrate to the true value would lead to overly pessimistic budgets to the point that the privacy system is not usable.

The data for Audience Engagement is being placed into Pinot on a daily basis and retained for 30 days. Hence, each day the data can potentially change, but not within the same day. We then use the query and the date to determine the pseudorandom seed, otherwise if we only use the query, then the noise added to a specific query would always be the same, despite the data changing. Note that we further use a secure key in the pseudorandom seed so that one cannot determine the seed only from the query and the date.

Recall that we are assuming analysts do not collude with each other and do not share the privatized results they receive, since this would mean essentially multiplying the information and call budgets that we enforced. However, with the pseudorandom seed being generated the same way for all analysts, we know that each analyst is receiving the same result so that even if they colluded, they could not average their results to get more confidence in the true result.

8.3. Rationale for Parameters in our Privacy System. There are ultimately four different parameters to set in our described privacy system: information budget k^* , call budget ℓ^* , ε_{per} , and δ . Note that with distinct counts, the ℓ_∞ -sensitivity is 1 for all queries and for any query that is deemed restricted sensitivity, we set $\Delta = 1$. The k parameter for top- k is an input from the analyst in each query and for unknown domain queries, \bar{d} is set to be as large as possible without harming the efficiency requirements of the product, typically $\bar{d} = \max\{10k, 1000\}$.

Part of the appeal of differential privacy is that it provides a worst case guarantee against privacy attacks, and can be simply stated as preventing an adversary from distinguishing whether a target's data was used in the analysis or not. This strong protection stops being very meaningful once the privacy loss parameter ε becomes large, say even larger than 1. However, deployments of differential privacy have quoted much larger privacy parameters than 1, see for example [4], and more recent works in private ML have used much larger parameters [25]. Further, these quoted parameters are for a one time calculation, rather than over multiple queries. Only deploying privacy systems that incorporate differential privacy with small ε would limit its applicability. In particular, in our setting we might only be able to allow for a single top-10 result before an analyst has exhausted his or her budget. Boiling down the entire privacy considerations of a complicated system to whether a couple of parameters stay below some arbitrary privacy threshold for all deployments seems overly simplistic. Other privacy safeguards can be added to increase the overall privacy, such as subsampling the dataset, as is done in this application. Further, the parameters in our system allow us to easily improve the overall privacy gains by modifying parameters. Providing the tuning knob between 100% utility and 100% privacy is incredibly helpful in showing the impact that differential privacy has on the overall product.

The question then is, what protections can differential privacy provide, even with large ε parameters. For this, we consider several different attacks, each used to set a certain parameter. This is not meant as an exhaustive list of all the attacks we considered, nor does it mean that these certain attacks are expected. This is merely to provide additional context to how privacy parameters can be set and might be useful for other privacy practitioners to use.

8.3.1. *Determining ε_{per} .* We consider the scenario where the dataset remains the same over the course of the data retention period (30 days) and each day an analyst asks the same query on that dataset. Recall how we set a pseudorandom seed for the same query and that it changes each day. Thus, the analyst would get 30 different noisy results on the same count. We then want to know the probability that the average of these noisy values will be within a tolerance of the true value. We set this tolerance to $1/2$, since this would mean that the analyst rounding to the nearest integer would reveal the true count. We then want to determine the following probability where $X_1, \dots, X_{30} \stackrel{i.i.d.}{\sim} \text{Lap}(2/\varepsilon_{\text{per}})$,

$$\Pr \left[\left| \frac{1}{30} \sum_{i=1}^{30} X_i \right| < 1/2 \right].$$

We approximate this probability with a Normal distribution, so that we have

$$\Pr \left[\left| \frac{1}{30} \sum_{i=1}^{30} X_i \right| < 1/2 \right] \approx \Pr \left[|\mathbb{N}(0, 1)| \leq \frac{\varepsilon_{\text{per}} \sqrt{30}}{4\sqrt{2}} \right].$$

Our aim is to reduce the chance of this attack, while also not adding too much noise to each count. We then sought to ensure roughly a 90% chance of no such attack. Plugging in $\varepsilon_{\text{per}} = 0.15$ leads to about an 11% probability of this attack being successful. Also recall that this attack will only be successful if the data does not change for the query over 30 days.

8.3.2. *Determining Information Budget.* One of the primary reasons for exploring differential privacy for this use case was differencing attacks. Consider the setting where an analyst asks two queries where they know that there is a single person different between the two in the unrestricted sensitivity setting. Despite the noise that we add to each count in each query, it is possible that the noise is small for some elements so it is clear what the true count was before noise. This happens when the noise is smaller than $1/2$ for a single count. Hence we want to compute the following probability, which can be written in terms of an exponential random variable Exp ,

$$p := \Pr[|\text{Lap}(2/\varepsilon_{\text{per}})| < 1/2] = \Pr[\text{Exp}(\varepsilon_{\text{per}}/2) < 1/2].$$

For a top- k query, we can expect to see pk (assume integer valued) many elements that have noisy count within half of the true count. If an adversary were to do a differencing attack, she would ask another top- k query and there would be fresh noise added. Hence, there would again be an expected pk many noisy counts that are close to the true counts. The adversary does not know which elements in both queries have noisy counts within $1/2$ of the true count, but we want to make sure these sets of elements do not overlap. If these elements with small noise do overlap in the two top- k results, then a difference between the two results might show the actual difference between the two.

We now want to prevent the possibility of these small count elements to overlap in the two top- k queries. Let's fix the set of pk elements that had noisy counts within half of the true counts in the first top- k . In the second top- k , we know that again there are expected to be pk elements, but where they are is random. Hence we get a uniformly random set of pk elements in the second top- k result and want to know what is the expected size of the

intersection between this set of pk elements and the pk elements from the first top- k . The probability that this intersection is of size s is the following:

$$\frac{\binom{pk}{s} \binom{k-pk}{pk-s}}{\binom{k}{pk}}.$$

This is a hypergeometric distribution and has expectation p^2k . To reduce the chance that these two sets of size pk overlap, we set the expected size of the intersection to be less than 1, i.e. $k < 1/p^2$.

Recall from the previous attack that we have $\varepsilon_{\text{per}} = 0.15$, we get $p = 0.0368$ and we then can use $k = 738$. Our information cost is the total number of results that can be returned plus one for the optimized threshold calculation in the unknown domain setting, which would bound the information cost by $2k + 2$, from these two top- k results. Note that we also return counts for the elements that we find, which also increases the total information cost by $2k$. Hence, we can set information budget $k^* = 4 \cdot 738 + 2 \approx 3000$, which from Figure 5, we see that more than 93% of analysts would not be impacted.

8.3.3. Determining δ and Call Budget. Our unknown domain algorithms include a threshold so that elements with a single user contribution (unique count) should not be shown in any result. However, noise is added to the threshold to ensure differential privacy, so we want to be able to control the chance that a unique count with noise becomes larger than the noisy threshold. Hence, we want to bound the probability that this can occur for a single count and then take a union bound over all \bar{d} counts that can be returned from Pinot. For the **UnkGumb** algorithm, we will write $Z_1, Z_2 \sim \text{Gumbel}(1/\varepsilon_{\text{per}})$ and $h_{(i)}$ as the i th ranked count in the input histogram \mathbf{h} . We then consider the following probability of a *bad single event* B_i where $i \leq \bar{d}$

$$\Pr[B_i] := \Pr[h_{(i)} + Z_1 > h_{(\bar{d}+1)} + 1 + \ln(\bar{d}/\delta)/\varepsilon_{\text{per}} + Z_2].$$

The worst case scenario is where every element in the histogram has count equal to 1, meaning only one user contributed to the counts and so each count is a unique count. Note that in **UnkGumb**, the threshold actually uses $\ln(\bar{k}/\delta)$ where there is an additional step to optimize from \bar{d} to a smaller \bar{k} , but here we use the larger \bar{d} to be more pessimistic. Noting that the difference $Z_1 - Z_2$ is distributed as a logistic random variable Log , we have the following

$$\begin{aligned} \Pr[B_i] &\leq \Pr \left[\text{Log}(1/\varepsilon_{\text{per}}) > 1 + \frac{\ln(\bar{d}/\delta)}{\varepsilon_{\text{per}}} \right] \\ &= 1 - \frac{1}{1 + e^{-\varepsilon_{\text{per}} \delta / \bar{d}}}. \end{aligned}$$

We then want to bound the event that any of the \bar{d} elements can appear above the threshold, hence

$$\begin{aligned} \Pr[\cup_{i=1}^{\bar{d}} B_i] &\leq \sum_{i=1}^{\bar{d}} \Pr[B_i] \\ &\leq \bar{d} \cdot \left(1 - \frac{1}{1 + e^{-\varepsilon_{\text{per}} \delta / \bar{d}}} \right) \leq \delta e^{-\varepsilon_{\text{per}}} \end{aligned}$$

Now, we want to make sure that the chance of this occurring over all queries is small. We use ℓ^* to denote the call cost which is how many top- k queries in the unknown domain setting are used and are the only types of queries that risk showing results that a single user contributed. Hence we will write B_i^ℓ to be a bad event for count $i \leq \bar{d}_\ell$ in the ℓ th query and take a union bound over all ℓ^* such queries.

$$\Pr[\cup_{\ell=1}^{\ell^*} \cup_{i=1}^{\bar{d}_\ell} B_i^\ell] \leq \ell^* \delta e^{-\varepsilon_{\text{per}}}.$$

We aim for a 1 in a hundred million chance and then also reference the cost budget usage in Figure 5. Setting $\ell^* = 30$ would not impact more than 95% of analysts and using $\delta = 10^{-10}$ with $\varepsilon_{\text{per}} = 0.15$ gives the overall probability bound of close to 1 in 400 million.

Use Case	Privacy Model	DP Algorithm Parameters (ϵ, δ)	Daily DP Parameters ($\epsilon_{\text{day}}, \delta_{\text{day}}$)	Monthly DP Parameters ($\epsilon_{\text{month}}, \delta_{\text{month}}$)
Google - RAPPOR [3] Chrome Homepages	Local ^b	(0.534, 0)	(25.63, 0) 30 min reporting	(769, 0)
Apple - Safari Domains [4]	Local	(4, 0)	(8, 0) ^a	(240, 0)
Apple - Emojis [4]	Local	(4, 0)	(4, 0) ^a	(120, 0)
Microsoft - Telemetry Collection per App [5]	Local ^b	(0.686, 0)	(2.74, 0) 6 hour reporting	(82.2, 0)
Google - Mobility Reports [26]	Global	(0.11, 0) or (0.22, 0)	(2.64, 0) ^c	(79.2, 0)
Microsoft - Assistive AI ^d	Global	(4, 10^{-7})	Not available	Not available
LinkedIn - Audience Engagement API ^e	Global	(0.15, 10^{-10})	—	(34.9, 7×10^{-9})

TABLE 2. Privacy parameters for existing deployments of privacy systems that use differentially private algorithms. Note that some parameters can be improved with a slightly larger δ_{month} .

^ahttps://www.apple.com/privacy/docs/Differential_Privacy_Overview.pdf

^bMemoization ensures that repeated records have much smaller overall privacy loss. This table shows the overall privacy loss for users that generate distinct records.

^cWe add up the parameters from [26] in daily visits in public places ($\epsilon = 1.74$), residential ($\epsilon = 0.44$), and workplaces ($\epsilon = 0.44$).

^d<https://www.microsoft.com/en-us/research/group/msai/articles/assistive-ai-makes-replying-easier-2>

^epre-GA status

8.4. Overall Privacy Guarantee. With these proposed parameter values for $\varepsilon_{\text{per}}, \delta, k^*, \ell^*$ and applying Theorem 2, we get a final $(34.9, 7 \times 10^{-9})$ -DP monthly guarantee. We want to put this guarantee in the context of other deployments of privacy systems that use differential privacy. Note that some privacy systems that adopted differential privacy are in the local privacy model, which is a more stringent privacy setting than the global model that we

consider here.³ Further, we used parameters that were publicly released, so they may differ from current deployments.

In Table 2, we identify different use cases that have adopted differential privacy and have released their privacy parameters along with how often data and reports are refreshed, thus allowing us to compute daily and monthly DP parameters. We focused on deployments where data is continually updated in both local and global models of privacy. It is important to point out that each privacy system includes additional safeguards beyond differentially private algorithms. Some of these differences include subsampling users, permuting records, and the use of memoization, as in Google's RAPPOR [3] and Microsoft's telemetry collection [5], to prevent longitudinal attacks when the same record is privatized with fresh noise repeatedly. What we account for in the table is if a user's data changes in the local model then fresh noise would be added to each result and hence the privacy loss accumulates.

9. CONCLUSION

We have presented a privacy system that incorporates state of the art algorithms for releasing histograms and top- k results in a differentially private way. Also, we have shown how we track the privacy budget for multiple analysts that can query our API. Combining the budget management service with DP algorithms allows us to make strong privacy guarantees of the overall system for any external partner that is allowed to make multiple, adaptively selected queries. This privacy system allows us to track the amount of information that is being released to external partners via the API in a precise way so that we can make informed decisions in how we can balance privacy safeguards with the usefulness of the product. We hope that this work demonstrates the feasibility of providing rigorous DP guarantees in systems that can scale.

Acknowledgements. We would like to thank Adrian Cardoso, Mark Cesar, Stephen Lynch, Sofus Macskassy, Koray Mancuhan, Sajjad Moradi, Sergey Yekhanin, and the entire LinkedIn Data Science Applied Research team for their helpful feedback on this work. Further, we thank Igor Perisic and Ya Xu for their support throughout this project.

³It is possible to compute the global DP parameters when local DP is used, see [27], [28], and [29]

REFERENCES

- [1] J.-F. Im, K. Gopalakrishna, S. Subramaniam, M. Shrivastava, A. Tumbde, X. Jiang, J. Dai, S. Lee, N. Pawar, J. Li, and R. Aringunram, “Pinot: Realtime olap for 530 million users,” in *Proceedings of the 2018 International Conference on Management of Data*, ser. SIGMOD ’18. New York, NY, USA: Association for Computing Machinery, 2018, p. 583–594. [Online]. Available: <https://doi.org/10.1145/3183713.3190661>
- [2] J. Dong, D. Durfee, and R. Rogers, “Optimal differential privacy composition for exponential mechanisms,” in *Proceedings of the 37th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, H. D. III and A. Singh, Eds., vol. 119. PMLR, 13–18 Jul 2020, pp. 2597–2606. [Online]. Available: <https://proceedings.mlr.press/v119/dong20a.html>
- [3] U. Erlingsson, V. Pihur, and A. Korolova, “Rappor: Randomized aggregatable privacy-preserving ordinal response,” in *Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS ’14. New York, NY, USA: ACM, 2014, pp. 1054–1067. [Online]. Available: <http://doi.acm.org/10.1145/2660267.2660348>
- [4] Apple Differential Privacy Team, “Learning with privacy at scale,” 2017, available at <https://machinelearning.apple.com/2017/12/06/learning-with-privacy-at-scale.html>.
- [5] B. Ding, J. Kulkarni, and S. Yekhanin, “Collecting telemetry data privately,” December 2017. [Online]. Available: <https://www.microsoft.com/en-us/research/publication/collecting-telemetry-data-privately/>
- [6] F. D. McSherry, “Privacy integrated queries: An extensible platform for privacy-preserving data analysis,” in *Proceedings of the 2009 ACM SIGMOD International Conference on Management of Data*, ser. SIGMOD ’09. New York, NY, USA: Association for Computing Machinery, 2009, p. 19–30. [Online]. Available: <https://doi.org/10.1145/1559845.1559850>
- [7] N. Johnson, J. P. Near, and D. Song, “Towards practical differential privacy for sql queries,” *Proc. VLDB Endow.*, vol. 11, no. 5, p. 526–539, Jan. 2018. [Online]. Available: <http://arxiv.org/abs/1706.09479>
- [8] K. Kenthapadi and T. T. L. Tran, “Pripearl: A framework for privacy-preserving analytics and reporting at linkedin,” in *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, ser. CIKM ’18. New York, NY, USA: Association for Computing Machinery, 2018, p. 2183–2191. [Online]. Available: <https://doi.org/10.1145/3269206.3272031>
- [9] M. Guevara, “Google developers,” Sep 2019. [Online]. Available: <https://developers.googleblog.com/2019/09/enabling-developers-and-organizations.html>
- [10] R. J. Wilson, C. Y. Zhang, W. Lam, D. Desfontaines, D. Simmons-Marengo, and B. Gipson, “Differentially private SQL with bounded user contribution,” *Proceedings on Privacy Enhancing Technologies*, vol. 2020, no. 2, pp. 230 – 250, 2020. [Online]. Available: <https://content.sciendo.com/view/journals/popets/2020/2/article-p230.xml>
- [11] A. N. Dajani, A. D. Lauger, P. E. Singer, D. Kifer, J. P. Reiter, A. Machanavajjhala, S. L. Garfinkel, S. A. Dahl, M. Graham, V. Karwa, H. Kim, P. Leclerc, I. M. Schmutte, W. N. Sexton, L. Vilhuber, and J. M. Abowd, “The modernization of statistical disclosure limitation at the U.S. Census bureau,” 2017, available online at <https://www2.census.gov/cac/sac/meetings/2017-09/statistical-disclosure-limitation.pdf>.
- [12] D. Durfee and R. M. Rogers, “Practical differentially private top-k selection with pay-what-you-get composition,” in *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. B. Fox, and R. Garnett, Eds., 2019, pp. 3527–3537. [Online]. Available: <https://proceedings.neurips.cc/paper/2019/hash/b139e104214a08ae3f2ebc149cdf6e-Abstract.html>
- [13] I. Kotsogiannis, Y. Tao, X. He, M. Fanaeepour, A. Machanavajjhala, M. Hay, and G. Miklau, “PrivateSQL: A differentially private SQL query engine,” *Proc. VLDB Endow.*, vol. 12, no. 11, p. 1371–1384, Jul. 2019. [Online]. Available: <https://doi.org/10.14778/3342263.3342274>
- [14] J. Kahan, “LinkedIn,” Sep 2019. [Online]. Available: <https://www.linkedin.com/pulse/microsoft-harvards-institute-quantitative-social-science-john-kahan/>
- [15] M. Gaboardi, J. Honaker, G. King, K. Nissim, J. Ullman, S. Vadhan, and J. Murtagh, “Psi (ψ): a private data sharing interface,” in *Theory and Practice of Differential Privacy*, New York, NY, 2016 2016. [Online]. Available: <https://arxiv.org/abs/1609.04340>
- [16] C. Dwork, F. McSherry, K. Nissim, and A. Smith, “Calibrating noise to sensitivity in private data analysis,” in *Proceedings of the Third Conference on Theory of Cryptography*, ser. TCC’06. Berlin, Heidelberg: Springer-Verlag, 2006, p. 265–284. [Online]. Available: https://doi.org/10.1007/11681878_14

- [17] C. Dwork, K. Kenthapadi, F. McSherry, I. Mironov, and M. Naor, "Our data, ourselves: Privacy via distributed noise generation," in *Advances in Cryptology - EUROCRYPT 2006*, S. Vaudenay, Ed. Berlin, Heidelberg: Springer Berlin Heidelberg, 2006, pp. 486–503. [Online]. Available: <https://www.iacr.org/archive/eurocrypt2006/40040493/40040493.pdf>
- [18] M. Cesar and R. Rogers, "Bounding, concentrating, and truncating: Unifying privacy loss composition for data analytics," in *Proceedings of the 32nd International Conference on Algorithmic Learning Theory*, ser. Proceedings of Machine Learning Research, V. Feldman, K. Ligett, and S. Sabato, Eds., vol. 132. PMLR, 16–19 Mar 2021, pp. 421–457. [Online]. Available: <https://proceedings.mlr.press/v132/cesar21a.html>
- [19] P. Kairouz, S. Oh, and P. Viswanath, "The composition theorem for differential privacy," vol. 63, no. 6, 2017, pp. 4037–4049. [Online]. Available: <https://doi.org/10.1109/TIT.2017.2685505>
- [20] J. Murtagh and S. Vadhan, "The complexity of computing the optimal composition of differential privacy," in *Proceedings, Part I, of the 13th International Conference on Theory of Cryptography - Volume 9562*, ser. TCC 2016-A. Berlin, Heidelberg: Springer-Verlag, 2016, pp. 157–175. [Online]. Available: https://doi.org/10.1007/978-3-662-49096-9_7
- [21] F. McSherry and K. Talwar, "Mechanism design via differential privacy," in *48th Annual IEEE Symposium on Foundations of Computer Science (FOCS'07)*, 2007, pp. 94–103. [Online]. Available: <http://dx.doi.org/10.1109/FOCS.2007.66>
- [22] A. Kafka, "A distributed streaming platform." [Online]. Available: kafka.apache.org/
- [23] L. Qiao, K. Surlaker, S. Das, T. Quiggle, B. Schulman, B. Ghosh, A. Curtis, O. Seeliger, Z. Zhang, A. Auradar, C. Beaver, G. Brandt, M. Gandhi, K. Gopalakrishna, W. Ip, S. Jgadesh, S. Lu, A. Pachev, A. Ramesh, A. Sebastian, R. Shanbhag, S. Subramaniam, Y. Sun, S. Topiwala, C. Tran, J. Westerman, and D. Zhang, "On brewing fresh espresso: LinkedIn's distributed data serving platform," in *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, ser. SIGMOD '13. New York, NY, USA: Association for Computing Machinery, 2013, p. 1135–1146. [Online]. Available: <https://doi.org/10.1145/2463676.2465298>
- [24] M. Bun and T. Steinke, "Concentrated differential privacy: Simplifications, extensions, and lower bounds," in *Theory of Cryptography*, M. Hirt and A. Smith, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2016, pp. 635–658. [Online]. Available: https://link.springer.com/chapter/10.1007/978-3-662-53641-4_24
- [25] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, "Deep learning with differential privacy," in *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, ser. CCS '16. New York, NY, USA: Association for Computing Machinery, 2016, p. 308–318. [Online]. Available: <https://doi.org/10.1145/2976749.2978318>
- [26] A. Aktay, S. Bavadekar, G. Cossoul, J. Davis, D. Desfontaines, A. Fabrikant, E. Gabrilovich, K. Gadepalli, B. Gipson, M. Guevara, C. Kamath, M. Kansal, A. Lange, C. Mandayam, A. Oplinger, C. Pluntke, T. Roessler, A. Schlosberg, T. Shekel, S. Vispute, M. Vu, G. Wellenius, B. Williams, and R. J. Wilson, "Google covid-19 community mobility reports: Anonymization process description (version 1.0)," 2020. [Online]. Available: <https://arxiv.org/abs/2004.04145>
- [27] U. Erlingsson, V. Feldman, I. Mironov, A. Raghunathan, K. Talwar, and A. Thakurta, "Amplification by shuffling: From local to central differential privacy via anonymity," in *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms*, ser. SODA '19. USA: Society for Industrial and Applied Mathematics, 2019, p. 2468–2479. [Online]. Available: <https://dl.acm.org/doi/10.5555/3310435.3310586>
- [28] B. Balle, J. Bell, A. Gascón, and K. Nissim, "The privacy blanket of the shuffle model," in *Advances in Cryptology - CRYPTO 2019 - 39th Annual International Cryptology Conference, Santa Barbara, CA, USA, August 18-22, 2019, Proceedings, Part II*, ser. Lecture Notes in Computer Science, A. Boldyreva and D. Micciancio, Eds., vol. 11693. Springer, 2019, pp. 638–667. [Online]. Available: https://doi.org/10.1007/978-3-030-26951-7_22
- [29] A. Cheu, A. D. Smith, J. Ullman, D. Zeber, and M. Zhilyaev, "Distributed differential privacy via shuffling," in *Advances in Cryptology - EUROCRYPT 2019 - 38th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Darmstadt, Germany, May 19-23, 2019, Proceedings, Part I*, ser. Lecture Notes in Computer Science, Y. Ishai and V. Rijmen, Eds., vol. 11476. Springer, 2019, pp. 375–403. [Online]. Available: https://doi.org/10.1007/978-3-030-17653-2_13

APPENDIX A. OMITTED ANALYSIS FOR SECTION 5.2

We now go through the analysis for Algorithm 3, in particular the proof of Lemma 5.3. The differences between Algorithm 3 and the version that appeared as Algorithm 4 in [12] is that we are returning counts as well as indices, we do not limit the number of outcomes to be at most k (since it is not a parameter), and we allow for counts to increase or decrease by $\tau \geq 1$ in neighboring datasets. As we will mainly be borrowing the analysis in [12] we will change \bar{d} to \bar{k} to better match the statements in that work. We then introduce the following algorithm, which we will show has the same distribution as $\text{UnkLap}^{\Delta, \bar{k}, \tau}(\mathbf{h})$.

Definition A.1 Limited Histogram Report Noisy Counts. *We assume that the ℓ_∞ sensitivity between any neighboring histograms is τ . We define the limited histogram report noisy counts to be $\text{LapMax}^{\bar{k}, \tau}$ that takes as input a histogram along with a domain set of indices and returns an ordered list of counts with the corresponding index, where $\text{LapMax}^{\bar{k}, \tau}(\mathbf{h}, \mathbf{d}) = (\{v_{(1)}, i_{(1)}\}, \dots, \{v_{(\perp)}, \perp\})$ and $(v_{(1)}, \dots, v_{(\perp)})$ is the sorted list of $v_i = h_{(i)} + \text{Lap}(2\tau\Delta/\varepsilon_{\text{per}})$ for each $i \in \mathbf{d}$ and $v_\perp = h_{(\bar{k}+1)} + \tau \left(1 + 2\Delta \ln(\Delta/\hat{\delta})/\varepsilon_{\text{per}}\right) + \text{Lap}(2\tau\Delta/\varepsilon_{\text{per}})$ with $\hat{\delta}$ given in Algorithm 3 as a function of $\delta > 0$.*

We have the following that connects $\text{LapMax}^{\bar{k}, \tau}$ with $\text{UnkLap}^{\Delta, \bar{k}, \tau}$.

Lemma A.1 . *For any histogram \mathbf{h} , we have that both mechanisms $\text{LapMax}^{\bar{k}, \tau}(\mathbf{h}, \mathbf{d}^{\bar{k}}(\mathbf{h}))$ and $\text{UnkLap}^{\bar{k}, \tau}(\mathbf{h})$ produce outcomes that are equal in distribution.*

If we fix a domain \mathbf{d} beforehand, then we have the following privacy statement. Note that the privacy of $\text{LapMax}^{\bar{k}, \tau}(\mathbf{h}, \mathbf{d})$ follows from the Laplace mechanism [16] being ε_{per} -DP. This is what allows us to output the counts as well as the indices. We just need to ensure that $i_{(\bar{k}+1)} \notin \mathbf{d}$ because if it was, then changing one index would change the count of both $h_{(\bar{k}+1)}$ and $h_\perp = h_{(\bar{k}+1)} + \tau(1 + \Delta \ln(\Delta/\delta)/\varepsilon_{\text{per}})$.

Lemma A.2 . *For any fixed $\mathbf{d} \subseteq [d]$ and neighbors \mathbf{h}, \mathbf{h}' such that $i_{(\bar{k}+1)}, i'_{(\bar{k}+1)} \notin \mathbf{d}$, then for any set of outcomes T ,*

$$\begin{aligned} \Pr[\text{LapMax}^{\bar{k}, \tau}(\mathbf{h}, \mathbf{d}) \in T] \\ \leq e^{\varepsilon_{\text{per}}/2} \Pr[\text{LapMax}^{\bar{k}, \tau}(\mathbf{h}', \mathbf{d}) \in T]. \end{aligned}$$

As was done in Durfee and Rogers [12], we can carefully account for the *good* (can bound the privacy loss) and *bad* (can bound these events with small probability) sets. Note that the outcome set of $\text{UnkLap}^{\Delta, \bar{k}, \tau}$ is a superset of Algorithm 4 in [12] when $k = \bar{k}$, and it is straightforward to see that these algorithms have the same distribution with respect to index output (ignoring the counts output from $\text{UnkLap}^{\Delta, \bar{k}, \tau}$). Therefore, all the bounds on the *bad* outcomes will still hold for our setting, and the analysis then follows from results in Section 6 of [12], where we state each result here.

Definition A.2 . *Given two neighboring histograms \mathbf{h}, \mathbf{h}' , we define \mathcal{S}_{Lap} as the outcome set of $\text{UnkLap}^{\Delta, \bar{k}, \tau}(\mathbf{h}, \mathbf{d}^{\bar{k}}(\mathbf{h}))$ (both indices and counts) and the outcome set of $\text{UnkLap}^{\Delta, \bar{k}, \tau}(\mathbf{h}', \mathbf{d}^{\bar{k}}(\mathbf{h}'))$ as $\mathcal{S}'_{\text{Lap}}$.*

We then define the bad outcomes as $\mathcal{S}_{\text{Lap}}^\delta := \mathcal{S}_{\text{Lap}} \setminus \mathcal{S}'_{\text{Lap}}$ and $\mathcal{S}'_{\text{Lap}}^\delta := \mathcal{S}'_{\text{Lap}} \setminus \mathcal{S}_{\text{Lap}}$.

Lemma A.3 . *For Δ -restricted sensitivity neighbors \mathbf{h}, \mathbf{h}' , we have*

$$\begin{aligned} & \Pr[\text{LapMax}^{\bar{k},\tau}(\mathbf{h}, \mathbf{d}^{\bar{k}}(\mathbf{h})) \in \mathcal{S}_{\text{Lap}}^\delta] \\ & \leq \delta/4 \cdot (3 + \ln(\Delta/\delta)) =: \bar{\delta} \end{aligned} \quad (\text{A.1})$$

Lemma A.4 . For any neighboring histograms \mathbf{h}, \mathbf{h}' and for any $S \subseteq \mathcal{S}_{\text{Lap}} \cap \mathcal{S}'_{\text{Lap}}$, we let $\mathbf{d}^{\varepsilon_{\text{per}}} = \mathbf{d}^{\bar{k}}(\mathbf{h}) \cap \mathbf{d}^{\bar{k}}(\mathbf{h}')$ and we must have the following for $\bar{\delta}$ given in (A.1)

$$\begin{aligned} & \Pr[\text{LapMax}^{\bar{k},\tau}(\mathbf{h}, \mathbf{d}^{\bar{k}}(\mathbf{h})) \in S] \\ & \leq \Pr[\text{LapMax}^{\bar{k},\tau}(\mathbf{h}, \mathbf{d}^{\varepsilon_{\text{per}}}) \in S] \\ & \leq \Pr[\text{LapMax}^{\bar{k},\tau}(\mathbf{h}, \mathbf{d}^{\bar{k}}(\mathbf{h})) \in S] + \bar{\delta} \end{aligned}$$

Lemma A.5 . For any neighboring histograms \mathbf{h}, \mathbf{h}' and any $S \subseteq \mathcal{S}_{\text{Lap}}$, then for $\bar{\delta}$ given in (A.1),

$$\begin{aligned} & \Pr[\text{LapMax}^{\bar{k},\tau}(\mathbf{h}, \mathbf{d}^{\bar{k}}(\mathbf{h})) \in S] \\ & \leq e^{\varepsilon_{\text{per}}/2} \Pr[\text{LapMax}^{\bar{k},\tau}(\mathbf{h}', \mathbf{d}^{\bar{k}}(\mathbf{h}')) \in S] \\ & \quad + (e^{\varepsilon_{\text{per}}/2} + 1)\bar{\delta}. \end{aligned}$$

Proof. We use the above results to get the following inequalities.

$$\begin{aligned} & \Pr[\text{LapMax}^{\bar{k},\tau}(\mathbf{h}, \mathbf{d}^{\bar{k}}(\mathbf{h})) \in S] \\ & = \Pr[\text{LapMax}^{\bar{k},\tau}(\mathbf{h}, \mathbf{d}^{\bar{k}}(\mathbf{h})) \in S \cap \{\mathcal{S}_{\text{Lap}} \cap \mathcal{S}'_{\text{Lap}}\}] \\ & \quad + \Pr[\text{LapMax}^{\bar{k},\tau}(\mathbf{h}, \mathbf{d}^{\bar{k}}(\mathbf{h})) \in S \cap \{\mathcal{S}_{\text{Lap}}^\delta\}] \\ & \leq \Pr[\text{LapMax}^{\bar{k},\tau}(\mathbf{h}, \mathbf{d}^{\varepsilon_{\text{per}}}) \in S \cap \{\mathcal{S}_{\text{Lap}} \cap \mathcal{S}'_{\text{Lap}}\}] \\ & \quad + \bar{\delta} \\ & \leq e^{\varepsilon_{\text{per}}/2} \\ & \quad \cdot \Pr[\text{LapMax}^{\bar{k},\tau}(\mathbf{h}', \mathbf{d}^{\varepsilon_{\text{per}}}) \in S \cap \{\mathcal{S}_{\text{Lap}} \cap \mathcal{S}'_{\text{Lap}}\}] + \bar{\delta} \\ & \leq e^{\varepsilon_{\text{per}}/2} \\ & \quad \left(\Pr[\text{LapMax}^{\bar{k},\tau}(\mathbf{h}', \mathbf{d}^{\bar{k}}(\mathbf{h}')) \in S \cap \{\mathcal{S}_{\text{Lap}} \cap \mathcal{S}'_{\text{Lap}}\}] + \bar{\delta} \right) \\ & \quad + \bar{\delta} \\ & \leq e^{\varepsilon_{\text{per}}/2} \Pr[\text{LapMax}^{\bar{k},\tau}(\mathbf{h}', \mathbf{d}^{\bar{k}}(\mathbf{h}')) \in S] + (e^{\varepsilon_{\text{per}}/2} + 1)\bar{\delta}. \end{aligned}$$

□