

# Private and Continual Release of Statistics

T-H. Hubert Chan  
The University of Hong Kong

Elaine Shi  
PARC

Dawn Song  
UC Berkeley

## Abstract

We ask the question – *how can websites and data aggregators continually release updated statistics, and meanwhile preserve each individual user’s privacy?* We propose a differentially private continual counter that outputs at every time step the approximate number of 1’s seen thus far. Our counter construction has error that is only poly-log in the number of time steps. We can extend the basic counter construction to allow websites to continually give top- $k$  and hot items suggestions while preserving users’ privacy.

## 1 Introduction

Websites such as online retailers, search engines and social networks commonly publish aggregate statistics about their users to realize valuable social and economic utilities. Moreover, the published statistics are continually updated over time as new data arrive. Such practices are ubiquitous and we name a few examples below. Sites such as Amazon, IMDB, Delicious and Flickr recommend popular items or content to users to enhance their browsing experience and engage their interests. Search engines such as Google and Yahoo help a user to auto-complete her search query by suggesting the most frequent search terms matching the prefix specified by the user. During political campaigns, websites survey the population and continually update the support rates for candidates.

Releasing aggregate information about users may seem harmless at first glance. However, previous work has shown that such statistical disclosures can expose sensitive information about an individual user [3, 12]. In particular, sites that continually update the published statistics over time can give even more leverage to the adversary and result in more severe privacy leakage [1].

In this paper, we ask the question – *how can we guarantee the users’ privacy when a website must continually publish new statistics as new data arrive?* Independent from our work, Dwork *et.al.* also consider essentially the same problem, and they phrase the problem as “*differential privacy under continual observation*” [7, 10, 11].

The setting we consider is different from the traditional setting in which differential privacy was studied. The traditional setting assumes a static database, and a curator who must answer  $k$  interactive queries or publish some sanitized statistics of the database non-interactively. In our setting, the database is dynamic and evolves over time, and a mechanism must update the published statistics as new data items arrive. Therefore, traditional differentially private mechanisms either fail to apply directly to our setting, or result in an unsatisfactory loss in terms of utility or privacy if applied naively.

### 1.1 Contributions

**Differentially private continual counter with poly-log error.** We consider the *continual counting problem*. Assume that the input stream  $\sigma \in \{0, 1\}^{\mathbb{N}}$  is a sequence of bits. The bit  $\sigma(t)$  at time  $t \in \mathbb{N}$  may denote whether an event of interest occurred at time  $t$ , e.g., whether a user purchased an item at time  $t$ . At every time step  $t \in \mathbb{N}$ , the mechanism must output an approximate count of the number of 1’s seen thus far.

We design an  $\epsilon$ -differentially private continual counter with small error. Specifically, for each  $t \in N$ , with probability at least  $1 - \delta$ , we guarantee  $O(\frac{1}{\epsilon} \cdot (\log t)^{1.5} \cdot \log \frac{1}{\delta})$  error<sup>1</sup>. This represents an improvement to Dwork’s recent

---

<sup>1</sup>For large values of  $t$ , we can actually get a better bound  $O(\frac{1}{\epsilon} \cdot (\log t)^{1.5} \cdot \sqrt{\log \frac{1}{\delta}})$ . To get a high probability statement, we can set  $\delta := \frac{1}{\text{poly}(t)}$  and the corresponding error becomes  $O((\log t)^2/\epsilon)$ .

result [7, 10] – they also design an  $\epsilon$ -differentially private continual counter, but with error  $O((\sqrt{n_t} + (\log t)^4)/\epsilon)$ , where  $n_t$  is the number 1’s in the stream up to time  $t$ .

Moreover, our mechanism achieves *time unboundedness*, i.e., the mechanism does not require a priori knowledge of an upper bound on the time for which it will run, and provides guarantees even when it is run indefinitely.

**Pan privacy.** Dwork *et.al.* first introduced the notion of pan privacy [7, 11]. A mechanism is pan privacy if it can preserve differential privacy even when an adversary can observe snapshots of the mechanism’s internal states, e.g., in subpoenas. We show how to modify our mechanism to achieve pan privacy, without incurring any loss in the asymptotic guarantees (Section 5).

**Applications.** Our continual counter construction has immediate practical applications. As mentioned earlier, it is a common practice for websites to suggest to users the most popular movies, news items or photos. In Section 6, we show how websites can continually make such top-k or hot items suggestions in a differentially private manner.

The counter is also an important primitive in numerous data streaming algorithms [2, 15, 17]. Our differentially private continual counter is an initial step towards designing a broad class of streaming algorithms that continually report outputs over time.

## 1.2 Related Work

**Most closely related work.** The most closely related work is an unpublished manuscript by Dwork *et.al.* [10] in which they construct a differentially private continual counter, with error square root in the number of 1’s seen thus far. Highlights of their results are published in a survey paper for an invited talk at the SODA’10 conference [7]. The results published in this paper are independent from their work, and represent an improvement over their results, as we achieve error that is independent of the number of 1’s in the stream and is only poly-log in the number of time steps.

Dwork, Naor, Pitassi and Rothblum [11] recently propose the notion of pan privacy, i.e., how to achieve differential privacy even in the presence of intrusions, in which the adversary is allowed access to the mechanism’s internal states. Dwork *et.al.* used the notion of pan privacy in the continual counter mechanism [7, 10], and showed how to make their counter mechanism resilient against a single unannounced intrusion. Inspired by their techniques, we also convert our mechanism to a pan private version that is immune to a single unannounced intrusion or multiple afterwards announced intrusions.

**Differential privacy in the traditional setting.** In the traditional setting, a trusted curator who holds a large data set must respond to queries *interactively* or publish sanitized statistics about the data *non-interactively*. The notion of differential privacy was first proposed and studied by Dwork *et.al.* [4, 9]. An extensive literature has since emerged, studying the different tradeoffs between utility and privacy. To better understand the motivation and state-of-the-art of this line of research, we recommend the readers to these excellent survey papers by Dwork [5, 6, 8].

Researchers have also applied theoretical results in differential privacy to real-world applications. For example, McSherry and Mironov show how to build privacy into the Netflix database published for the Netflix contest [16]. Korolova *et.al.* show how to release search logs and click logs privately [14].

**Attacks against privacy.** A complementary line of research is attacks against privacy. Narayanan *et.al.* show how to de-anonymize the Netflix data set [18]. Jones *et.al.* show how to break the privacy of query log bundles [13]. More relevant to this work, Calandrino *et.al.* [1] recently demonstrate that by observing continual updates from websites such as Amazon over a period of time, an adversary can learn individual user behavior at a fine level of granularity. Our work is partly inspired by the problem they expose.

## 2 Preliminaries

### 2.1 Continual Counting Mechanism

We consider streams of 0's and 1's. Formally, a stream  $\sigma \in \{0, 1\}^{\mathbb{N}}$  is a bit-string of countable length, where  $\mathbb{N} := \{1, 2, 3, \dots\}$  is the set of positive integers. Specifically,  $\sigma(t) \in \{0, 1\}$  denotes the bit at time  $t \in \mathbb{N}$ . We write  $[T] := \{1, 2, 3, \dots, T\}$  and  $\sigma_T \in \{0, 1\}^T$  is the length  $T$  prefix of the stream  $\sigma$ . We will use the term *item* to refer to a bit in the stream.

At every time  $t$ , we wish to output the number of 1's that have arrived up to time  $t$ .

**Definition 1** (Continual Counting Query). *Given a stream  $\sigma \in \{0, 1\}^{\mathbb{N}}$ , the count for the stream is a mapping  $c_\sigma : \mathbb{N} \rightarrow \mathbb{Z}$  such that for each  $t \in \mathbb{N}$ ,  $c_\sigma(t) := \sum_{i=1}^t \sigma(i)$ . We write  $c$  instead of  $c_\sigma$  when there is no risk of ambiguity on the stream  $\sigma$  in question.*

We now formally define the notion of a continual counting mechanism which continually outputs the number of 1's seen thus far.

**Definition 2** (Counting Mechanism). *A counting mechanism  $\mathcal{M}$  takes a stream  $\sigma \in \{0, 1\}^{\mathbb{N}}$  and produces a (possibly randomized) mapping  $\mathcal{M}(\sigma) : \mathbb{N} \rightarrow \mathbb{R}$ . Moreover, for all  $t \in \mathbb{N}$ ,  $\mathcal{M}(\sigma)(t)$  is independent of all  $\sigma(i)$ 's for  $i > t$ . We can also view  $\mathcal{M}(\sigma)$  as a point in  $\mathbb{R}^{\mathbb{N}}$ . When there is no risk of ambiguity on the stream  $\sigma$  in question, we drop the dependence on  $\sigma$  and use  $\mathcal{M}(t)$  to mean  $\mathcal{M}(\sigma)(t)$ .*

**Definition 3** (Time-bounded Mechanism). *A counting mechanism  $\mathcal{M}$  is unbounded, if it accepts streams of indefinite lengths, i.e., given any stream  $\sigma$ ,  $\mathcal{M}(\sigma) \in \mathbb{R}^{\mathbb{N}}$ . Given  $T \in \mathbb{N}$ , a mechanism  $\mathcal{M}$  is  $T$ -bounded if it only accepts streams of lengths at most  $T$  and returns  $\mathcal{M}(\sigma) \in \mathbb{R}^T$ . In other words, the mechanism needs to know the value  $T$  in advance and only looks at the length  $T$  prefix of any given stream.*

We would like the mechanism to be useful, that is, its output should well approximate the true count at any point of time. We formally define the notion of utility below.

**Definition 4** (Utility). *A counting mechanism  $\mathcal{M}$  is  $(\lambda, \delta)$ -useful at time  $t$ , if for any stream  $\sigma$ , with probability at least  $1 - \delta$ , we have  $|c_\sigma(t) - \mathcal{M}(\sigma)(t)| \leq \lambda$ . Note that  $\lambda$  may be a function of  $\delta$  and  $t$ .*

### 2.2 Differential Privacy

Now consider the notion of differential privacy in this setting. Intuitively, a mechanism is differentially private if its outputs on any two “adjacent” streams are roughly the same. This means that an adversary is unable to determine whether an event of interest took place or not by observing the output of the mechanism over time. For example, the adversary is unable to determine whether a user purchased an item at some time  $t$ .

**Definition 5** (Differential Privacy). *Two streams  $\sigma$  and  $\sigma'$  are adjacent if they differ at exactly one time  $t$ . A counting mechanism  $\mathcal{M}$  is  $\epsilon$ -differentially private (or preserves  $\epsilon$ -differential privacy) if for any adjacent streams  $\sigma$  and  $\sigma'$ , and any measurable subset  $S \subseteq \mathbb{R}^{\mathbb{N}}$  (or  $S \subseteq \mathbb{R}^T$  for  $T$ -bounded mechanisms),  $\Pr[\mathcal{M}(\sigma) \in S] \leq \exp(\epsilon) \cdot \Pr[\mathcal{M}(\sigma') \in S]$ .*

Dwork *et al.* first formalized the notion of pan privacy [7, 11] to deal with intruders who can observe snapshots of the mechanism's internal states, e.g., in a subpoena. We have to assume some notion of “atomicity”, that is, an intrusion can happen only when the mechanism has finished its update at a certain time step. During the update, the mechanism has to store the true value of new item  $\sigma(t)$  somehow, and intrusion at this point would clearly break the privacy.

**Definition 6** (Pan Privacy against Single Unannounced Intrusion). *Suppose  $I$  is the set of the internal states of a mechanism  $\mathcal{M}$ . Given a stream  $\sigma$  and an unannounced intrusion at time  $t$ , we can view the output  $\mathcal{M}(\sigma, t)$  of the mechanism as an element  $(i, s) \in I \times \mathbb{R}^{\mathbb{N}}$ , where  $i$  represents the knowledge gained in the one single intrusion at time  $t$ . A counting mechanism  $\mathcal{M}$  is  $\epsilon$ -pan private (or preserves  $\epsilon$ -pan privacy) against single intrusion if for any adjacent streams  $\sigma$  and  $\sigma'$ , any time  $t$  and any measurable subset  $S \subseteq I \times \mathbb{R}^{\mathbb{N}}$ ,  $\Pr[\mathcal{M}(\sigma, t) \in S] \leq \exp(\epsilon) \cdot \Pr[\mathcal{M}(\sigma', t) \in S]$ .*

**Definition 7** (Pan Privacy against Multiple Afterwards Announced Intrusions). *Suppose  $I$  is the set of the internal states of a mechanism  $\mathcal{M}$ . Let  $K \subset \mathbb{N}$  be a subset of size  $k$  that represents the time steps at which intrusions are made. We assume that an intrusion is not known in advance, but the mechanism is aware of an intrusion immediately after it has happened. Given a stream  $\sigma$  and the intrusion times  $K$ , we can view the output  $\mathcal{M}(\sigma, K)$  as an element  $(i, s) \in I^k \times \mathbb{R}^{\mathbb{N}}$ , where  $i$  represents the knowledge gained in the  $k$  intrusions. A counting mechanism  $\mathcal{M}$  is  $\epsilon$ -pan private (or preserves  $\epsilon$ -pan privacy) against multiple intrusions if for any adjacent streams  $\sigma$  and  $\sigma'$ , any  $K \subset \mathbb{N}$ , and any measurable subset  $S \subseteq I^{|K|} \times \mathbb{R}^{\mathbb{N}}$ ,  $\Pr[\mathcal{M}(\sigma, K) \in S] \leq \exp(\epsilon) \cdot \Pr[\mathcal{M}(\sigma', K) \in S]$ .*

## 2.3 Tools

In the design of differentially private mechanisms, the Laplace distribution is often used to introduce random noise [4, 9]. We use  $\text{Lap}(b)$  to denote the Laplace distribution with mean 0 and variance  $2b^2$ . Its probability density function is  $x \mapsto \frac{1}{2b} \exp(-\frac{|x|}{b})$ .

Dwork *et.al.* showed that if we mask the true answer of a query with Laplacian noise proportional to the sensitivity of the query function, such a mechanism preserves differential privacy for static databases [4, 9]. This is stated formally in the Fact 1.

**Fact 1** (Laplace Distribution Maintains Differential Privacy.). *Let  $a, b \in \mathbb{R}$  and  $|a - b| \leq \Delta$ . Let  $\gamma \sim \text{Lap}(\frac{\Delta}{\epsilon})$  be a random variable having Laplace distribution. Then, for any measurable subset  $S \subseteq \mathbb{R}$ ,  $\Pr[a + \gamma \in S] \leq \exp(\epsilon) \cdot \Pr[b + \gamma \in S]$ .*

In the constructions that we propose, the noise may not come from a single Laplace distribution, but rather is the sum of multiple independent Laplace distributions. We now derive a property of the sum of independent Laplace distributions.

**Lemma 1** (Sum of Independent Laplace Distributions). *Suppose  $\gamma_i$ 's are independent random variables, where each  $\gamma_i$  has Laplace distribution  $\text{Lap}(b_i)$ . Suppose  $Y := \sum_i \gamma_i$ , and  $b_M := \max_i b_i$ . Let  $\nu \geq \sqrt{\sum_i b_i^2}$  and  $0 < \lambda < \frac{2\nu^2}{b_M}$ . Then,  $\Pr[Y > \lambda] \leq \exp(-\frac{\lambda^2}{8\nu^2})$ .*

*Proof.* We use moment generating functions in a Chernoff-like argument. For each  $\gamma_i$ , the moment generating function is

$$E[\exp(h\gamma_i)] = \frac{1}{1-h^2b_i^2}, \text{ where } |h| < \frac{1}{b_i}.$$

Using the inequality  $(1-x)^{-1} \leq 1+2x \leq \exp(2x)$ , for  $|x| < \frac{1}{2}$ , we have  $E[\exp(h\gamma_i)] \leq \exp(2h^2b_i^2)$ , if  $|h| < \frac{1}{2b_i}$ .

We next use a standard calculation. For  $0 < h < \frac{1}{b_M}$ , we have

$$\begin{aligned} \Pr[Y > \lambda] &= \Pr[\exp(hY) > \exp(h\lambda)] \\ &\leq \exp(-h\lambda) E[\exp(hY)] \\ &= \exp(-h\lambda) \prod_i E[\exp(h\gamma_i)] \\ &\leq \exp(-h\lambda + 2h^2\nu^2) \end{aligned}$$

By assumption,  $0 < \lambda < \frac{2\nu^2}{b_M}$ . Setting  $h := \frac{\lambda}{4\nu^2} < \frac{1}{2b_M}$ , we conclude that  $\Pr[Y > \lambda] \leq \exp(-\frac{\lambda^2}{8\nu^2})$ .  $\square$

**Corollary 1** (Measure Concentration). *Let  $Y$ ,  $\nu$ ,  $\{b_i\}_i$  and  $b_M$  be defined as in Lemma 1. Suppose  $0 < \delta < 1$  and  $\nu > \max\{\sqrt{\sum_i b_i^2}, b_M \sqrt{2 \ln \frac{2}{\delta}}\}$ . Then,  $\Pr[|Y| > \nu \sqrt{8 \ln \frac{2}{\delta}}] \leq \delta$ .*

*To simplify our presentation and improve readability, we choose  $\nu := \sqrt{\sum_i b_i^2} \cdot \sqrt{2 \ln \frac{2}{\delta}}$  and use the following slightly weaker result: with probability at least  $1 - \delta$ , the quantity  $|Y|$  is at most  $O(\sqrt{\sum_i b_i^2} \log \frac{1}{\delta})$ .*

### 3 Time-Bounded Counting Mechanisms

In this section, we describe mechanisms that require a priori knowledge of an upper bound on time. In Section 4.2, we show how to remove this requirement, and achieve unbounded counting mechanisms.

#### 3.1 Simple Counting Mechanisms

To aid the understanding of our contributions and techniques, we first explain two simple constructions.

**Simple Counting Mechanism I.** The mechanism is given a stream  $\sigma \in \{0, 1\}^{\mathbb{N}}$ , a differential privacy parameter  $\epsilon > 0$ , and an upper bound  $T$  on time. At each time step  $t$ , the mechanism samples a fresh independent random variable  $\gamma_t \sim \text{Lap}(\frac{1}{\epsilon})$ , and releases  $\alpha_t = c(t) + \gamma_t$ , where  $c(t)$  is the true count at time step  $t$ . It is not hard to see that the above mechanism is  $O(T\epsilon)$ -differentially private, and at each time step, the error is  $O(\frac{1}{\epsilon})$  with high probability. Alternatively, one can substitute  $\epsilon' = \epsilon/T$ , and add much bigger noise  $\sim \text{Lap}(\frac{1}{\epsilon'})$  at every time step. In this way, we get  $\epsilon$  differential privacy; however, now the error at each time step is  $O(\frac{T}{\epsilon})$ .

Simple mechanism I is a straightforward extension of the Laplace mechanism proposed by Dwork *et.al.* [4, 9]. Basically, at every time step, the mechanism answers a new query, and randomizes the answer with fresh independent noise. The down side of this approach is that the privacy loss grows linearly with respect to the number of queries, which is  $t$  in our setting.

**Simple Counting Mechanism II.** In essence, Simple Counting Mechanism II produces a “sanitized” stream by adding independent Laplacian noise to each item in the stream. Suppose the mechanism is given a stream  $\sigma \in \{0, 1\}^{\mathbb{N}}$  and a differential privacy parameter  $\epsilon > 0$ . For each time step  $t \in \mathbb{N}$ , the mechanism samples an independent random variable  $\gamma_t$  with Laplace distribution  $\text{Lap}(\frac{1}{\epsilon})$ . Define  $\alpha_t := \sigma(t) + \gamma_t$ . Then, the mechanism  $\mathcal{M}$  gives the output  $\mathcal{M}(\sigma)(t) := \sum_{i \leq t} \alpha_i$  at time  $t$ . A similar idea has been proposed as a survey technique by Warner [19].

It is not hard to see that Simple Mechanism II can be implemented with  $O(1)$  words of memory and is unbounded and  $\epsilon$ -differentially private. We use Corollary 1 to analyze the utility of the mechanism. Fix some time  $T$ . Observe that  $\mathcal{M}(\sigma)(T) - c_\sigma(T) = \sum_{t \leq T} \gamma_t =: Y$ . In this case, all  $\gamma_t \sim \text{Lap}(\frac{1}{\epsilon})$ . Hence, all  $b_t := \frac{1}{\epsilon}$ .

**Theorem 1.** *Let  $0 < \delta < 1$ ,  $\epsilon > 0$ . the Simple Counting Mechanism II is  $\epsilon$ -differentially private, and is  $(O(\frac{\sqrt{t}}{\epsilon} \cdot \log \frac{1}{\delta}), \delta)$ -useful at any time  $t \in \mathbb{N}$ .*

#### 3.2 Intuition

We will describe the Two-Level Counting Mechanism and the Binary Counting Mechanism. Informally, the Two-Level Mechanism achieves  $\epsilon$ -differential privacy and  $O(t^{\frac{1}{4}})$  error. The Binary Mechanism is a further improvement, and achieves  $O((\log t)^{1.5})$  error while maintaining  $\epsilon$ -differential privacy. We now explain the intuitions for the Two-Level Mechanism and the Binary Mechanism.

**A framework for describing mechanisms.** We will describe our counting mechanisms using a common framework. Recall that the job of the mechanism is to output an approximate count at every time. However, from now on, we will think of our mechanisms as releasing noisy “p-sums” instead of counts. One can think of p-sums as intermediate results from which an observer can estimate the count at every time step herself.

**Definition 8 (p-sum).** *A p-sum is a partial sum of consecutive items. Let  $1 \leq i \leq j$ . We use the notation  $\Sigma[i, j] := \sum_{k=i}^j \sigma(k)$  to denote a partial sum involving items  $i$  through  $j$ .*

Furthermore, once we add noise to a p-sum, we obtain a noisy p-sum denoted as  $\widehat{\Sigma}$ .

The mechanisms we consider will release noisy versions of these p-sums as new items arrive. When an observer sees the sequence of p-sums, she can compute an estimate for the count at each time step, in particular, by summing up an appropriate selection of p-sums. For example, if an observer sees a noisy p-sum  $\widehat{\Sigma}[1, k] = \Sigma[1, k] + \text{noise}$  released at time step  $k$ , and another noisy p-sum  $\widehat{\Sigma}[k+1, t] = \Sigma[k+1, t] + \text{noise}$  released at time step  $t$ , then she

can estimate the count at time  $t$  by summing up these two noisy  $\mathbf{p}$ -sums, i.e.,  $\widehat{\Sigma}[1, k] + \widehat{\Sigma}[k + 1, t]$ . Notice that the observer needs to be able to do this not only for a specific time  $t$ , but also for every time step in  $\mathbb{N}$ .

Now we rethink Simple Mechanism I using this framework. The noisy  $\mathbf{p}$ -sums released are noisy versions of the true count for each time step, that is,  $\{\widehat{\Sigma}[1, t] = \Sigma[1, t] + \text{noise}\}_{1 \leq t \leq T}$ , where  $\Sigma[1, t] = c(t)$  is the true count at time  $t$ . In this case, the  $\widehat{\Sigma}[1, t]$  itself is the estimated count at time  $t$ ; and therefore can be regarded as a sum of noisy  $\mathbf{p}$ -sums (with only one summand). Notice that each item  $\sigma(t)$  appears in  $O(T)$  of these  $\mathbf{p}$ -sums. This means that when you flip an item in the incoming stream,  $O(T)$  of these  $\mathbf{p}$ -sums will be affected – this is the reason why the privacy loss is linear in  $T$ .

Now consider Simple Mechanism II. The noisy  $\mathbf{p}$ -sums released are noisy versions of each item  $\widehat{\Sigma}_t = \Sigma[t, t] + \text{noise}$ , where  $\Sigma[t, t] = \sigma(t)$  is the  $t$ -th item itself. In this case, each item appears in only one  $\mathbf{p}$ -sum, however, each count is the sum of  $O(T)$   $\mathbf{p}$ -sums. More specifically, to estimate the count at time  $t$ , the observer sums up  $t$  noisy  $\mathbf{p}$ -sums  $\widehat{\Sigma}_1, \dots, \widehat{\Sigma}_t$ . As each noisy  $\mathbf{p}$ -sum contains some fresh independent noise, the noises add up. In fact, over  $t$  time steps, the error would be  $O(\sqrt{t})$  with high probability.

**Observation 1.** (Informal.) Suppose a mechanism  $\mathcal{M}$  adds  $\text{Lap}(\frac{1}{\epsilon})$  noise to every  $\mathbf{p}$ -sum before releasing it. In  $\mathcal{M}$ , each item in the stream appears in at most  $x$   $\mathbf{p}$ -sums, and each estimated count is the sum of at most  $y$   $\mathbf{p}$ -sums. Then, the mechanism  $\mathcal{M}$  achieves  $x \cdot \epsilon$  differential privacy. Moreover, from Corollary 1, the error is  $O(\frac{\sqrt{y}}{\epsilon})$  with high probability. Alternatively, to achieve  $\epsilon$ -differential privacy, one can scale appropriately by having  $\epsilon' = \frac{\epsilon}{x}$ . Now if the mechanism instead adds  $\text{Lap}(\frac{1}{\epsilon'})$  noise to each  $\mathbf{p}$ -sum, we achieve  $\epsilon$ -differential privacy, and  $O(\frac{x\sqrt{y}}{\epsilon})$  error with high probability.

Mechanism	Each item appears in ? $\mathbf{p}$ -sums	Each count is the sum of ? $\mathbf{p}$ -sums	Asymptotic error (while maintaining $\epsilon$ diff. priv.)
Simple I	$O(T)$	$O(1)$	$O(T)$
Simple II	$O(1)$	$O(T)$	$O(\sqrt{T})$
Two-Level	$O(1)$	$O(\sqrt{T})$	$O(T^{\frac{1}{4}})$
Binary	$O(\log T)$	$O(\log T)$	$O((\log T)^{1.5})$

Table 1: Informal intuition for the Two-Level Mechanism and the Binary Mechanism. For simplicity, we omit the parameters  $\epsilon$  and  $\delta$  from the bounds.

**Goal.** From the above analysis, it appears that an inherent tension exists between utility (i.e., small error) and privacy, and our challenge is how to strike a balance between the two conflicting goals. We would like to achieve the following goals.

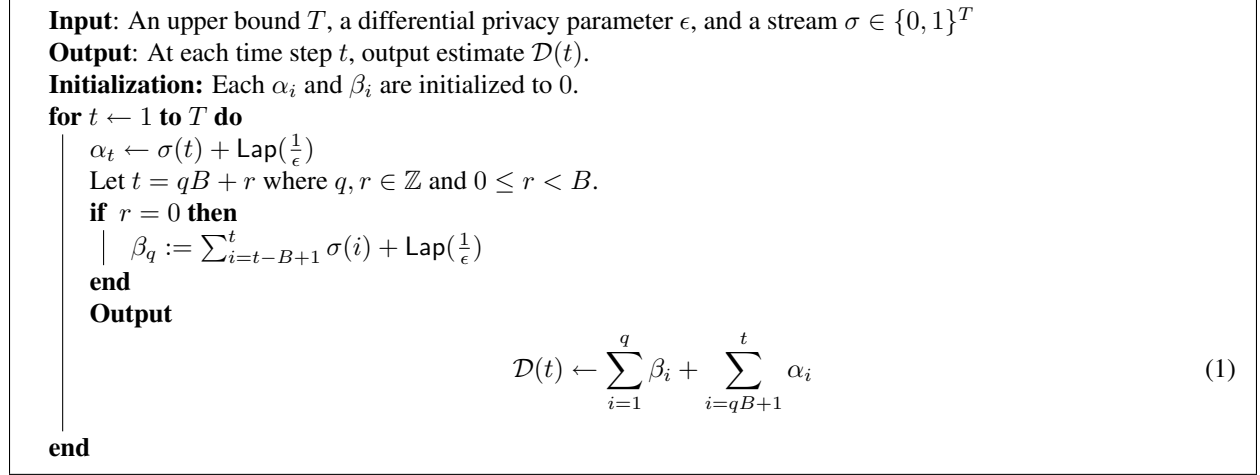
- *Each item appears in a small number of  $\mathbf{p}$ -sums*. Intuitively, this limits the influence of any item and guarantees small privacy loss. More specifically, when one flips an item in the incoming stream, not too many  $\mathbf{p}$ -sums will be affected.
- *Each count is a sum of a small number of  $\mathbf{p}$ -sums*. Each noisy  $\mathbf{p}$ -sum contains some noise, and the noises add up as one sums up several noisy  $\mathbf{p}$ -sums. However, if each output count is the sum of a small number of noisy  $\mathbf{p}$ -sums, the accumulation of noises is bounded. In this way, we can achieve small error.

Table 1 summarizes how the Two-Level mechanism and the Binary Mechanism utilize the above ideas to achieve small error while preserving strong privacy.

### 3.3 Two-Level Counting Mechanism

Using the Simple Counting Mechanism II as a building block, we describe the Two-Level Counting Mechanism. The idea is that when items from the stream come, we group them in contiguous blocks of size  $B$ . Within a block, we

run the Simple Counting Mechanism II. On top of that, we run another Simple Counting Mechanism II, treating each block as a single element.



**Figure 1:** Two-Level Mechanism  $\mathcal{D}$

**Two-Level Mechanism: the p-sum view.** One good way to understand the Two-Level Mechanism is to consider it under the p-sum framework. Notice that in Figure 1, each  $\beta_q = \widehat{\Sigma}[(q-1)B+1, qB]$  is a noisy p-sum, and each  $\alpha_t = \widehat{\Sigma}[t, t]$  is also a noisy p-sum (for a single item). It would suffice if the mechanism simply released the set of noisy p-sums:

$$\{\beta_q | 1 \leq q \leq \lfloor T/B \rfloor\} \cup \{\alpha_t | 1 \leq t \leq T\},$$

as an observer can reconstruct the approximate count at any time step from these noisy p-sums, according to Equation (1) of Figure 1.

Observe that each item  $\sigma(t)$  appears in at most two p-sums: at most one of the  $\beta$ 's and at most one of the  $\alpha$ 's. Specifically, let  $q := \lceil \frac{t}{B} \rceil$ , then  $\sigma(t)$  appears in only  $\beta_q$  and  $\alpha_t$ . Hence, we can conclude that the counting mechanism preserves  $2\epsilon$ -differential privacy.

From Equation (1), it is not hard to see that the estimated count at any time  $t$  is the sum of at most  $\lfloor t/B \rfloor + B$  noisy p-sums. In particular, if we let  $B = \sqrt{T}$ , then the estimated count at any time is the sum of at most  $2B$  noisy p-sums. According to Observation 1, the error is roughly  $O(T^{1/4}/\epsilon)$  with high probability.

We formalize the above intuition with the following theorem.

**Theorem 2.** *The Two-Level Counting Mechanism is  $2\epsilon$ -differentially private. Furthermore, for each  $t \in \mathbb{N}$ , the Two-Level Counting Mechanism with block size  $B$  is  $(O(\frac{1}{\epsilon} \cdot \sqrt{(\frac{t}{B} + B) \cdot \log \frac{1}{\delta}}), \delta)$ -useful at time  $t$ .*

*Proof.* The differential privacy argument is straightforward. As mentioned above, each item in the stream  $\sigma(t)$  appears in at most 2 noisy p-sums. Therefore, if we flip  $\sigma(t)$ , at most 2 noisy p-sums will be affected.

We now prove the utility part of the theorem. Observe that at any time  $t = qB + r$  where  $q, r \in \mathbb{Z}$  and  $0 \leq r < B$ , the error  $\mathcal{D}(t) - c_\sigma(t)$  is the sum of  $K = q + r$  independent Laplacian distributions  $\text{Lap}(\frac{1}{\epsilon})$ . Since  $\frac{t}{B} \leq K \leq (\frac{t}{B} + B)$ , it follows as in the proof for the Simple Counting Mechanism II that, at time  $T$ , the Two-Level Counting Mechanism is  $(O(\frac{1}{\epsilon} \cdot \sqrt{(\frac{t}{B} + B) \cdot \log \frac{1}{\delta}}), \delta)$ -useful at time  $t$ .  $\square$

Given  $T \in \mathbb{N}$ , we can set  $B := \lfloor \sqrt{T} \rfloor$  to form a  $T$ -bounded counting mechanism.

**Corollary 2.** *Let  $0 < \delta < 1$  and  $\epsilon > 0$ . For each  $T \in \mathbb{N}$ , there is a  $T$ -bounded counting mechanism that preserves  $2\epsilon$ -differential privacy and is  $(O(\frac{1}{\epsilon} \cdot T^{1/4} \cdot \log \frac{1}{\delta}), \delta)$ -useful at each time  $t \in [T]$ .*

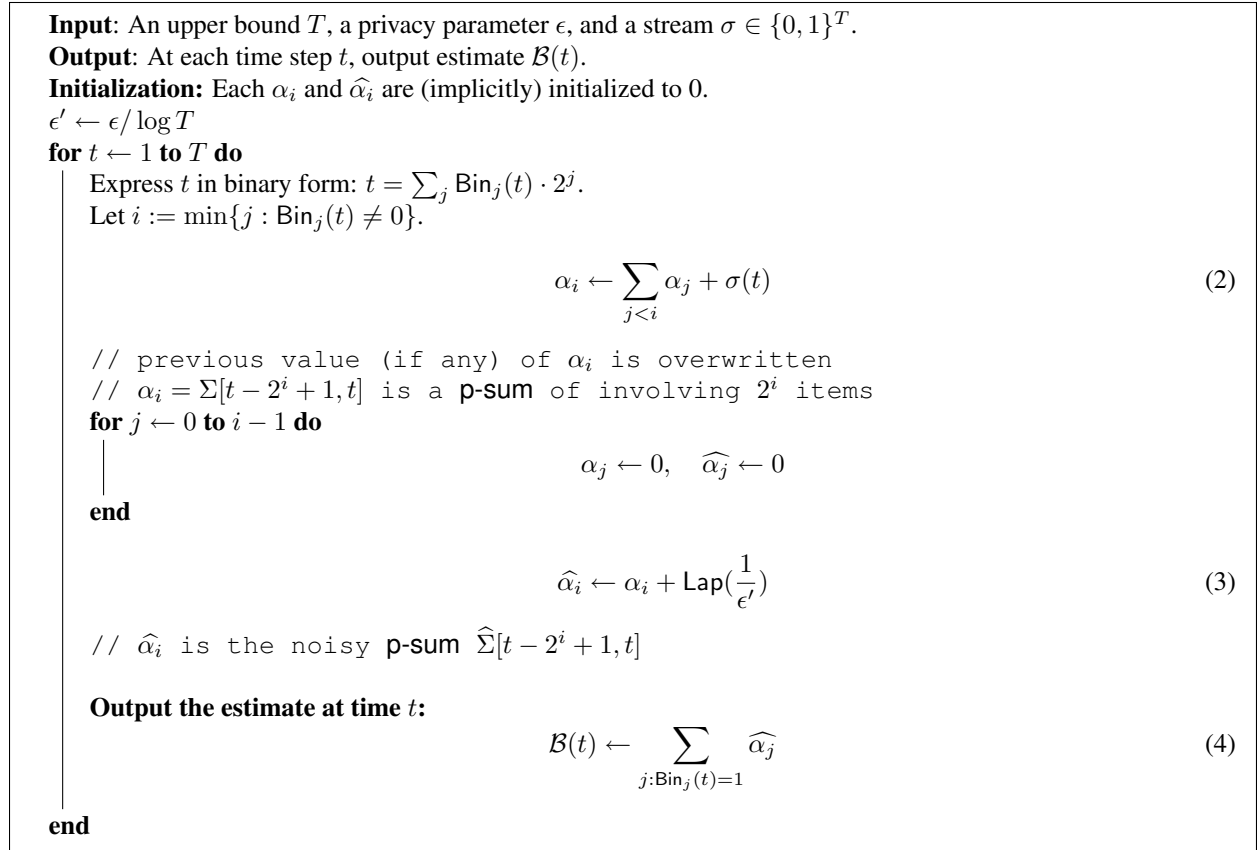
Finally, we point out that the Two-Level Mechanism can actually be implemented with  $O(1)$  memory. At any time  $t = qB + r$  where  $q, r \in \mathbb{Z}$  and  $0 \leq r < B$ , the mechanism only needs to store the following values:  $\sum_{i=1}^q \beta_i$ ,  $\Sigma[qB+1, t]$  and  $\sum_{i=qB+1}^t \alpha_i$ . Hence, we have the following claim.

**Claim 1.** *The Two-Level Counting Mechanism can be implemented with  $O(1)$  words of memory.*

### 3.4 Binary Counting Mechanism

We could extend the idea of the Two-Level Counting Mechanism to a Multi-level Counting Mechanism, and compute the optimal number of levels given  $T$ , the upper bound on time. However, we take a better approach called the Binary Mechanism. The idea is that at any time  $t$ , the counting mechanism internally groups the items that have arrived to form  $\mathbf{p}$ -sums of different sizes. The precise grouping of the items depends on the binary representation of the number  $t$  – hence the name Binary Mechanism

Given any number  $t \in N$ , let  $\text{Bin}_i(t) \in \{0, 1\}$  be the  $i$ th digit in the binary representation of  $t$ , where  $\text{Bin}_0(t)$  is the least significant digit. Hence,  $t = \sum_i \text{Bin}_i(t) \cdot 2^i$ . Informally, if  $\text{Bin}_i(t) = 1$ , then there is a  $\mathbf{p}$ -sum involving  $2^i$  items. We formally describe the Binary Mechanism in Figure 2.



**Figure 2:** Binary Mechanism  $\mathcal{B}$

**Binary mechanism: the  $\mathbf{p}$ -sum view.** The best way to understand the Binary Mechanism is to think in terms of the  $\mathbf{p}$ -sum framework described earlier. Basically, instead of outputting the estimated counts, the mechanism could equivalently release a sequence of noisy  $\mathbf{p}$ -sums which provide sufficient information for an observer to estimate the count at each time step  $t$ . In particular, at any time  $t$ , the Binary Mechanism releases a new noisy  $\mathbf{p}$ -sum  $\hat{\Sigma}[t - 2^i + 1, t]$  of length  $2^i$  and ending at position  $t$ , where  $i$  denotes the position of the least significant non-zero bit in the binary representation of  $t$ . This  $\mathbf{p}$ -sum and its noisy version are (temporarily) saved in the variables  $\alpha_i$  and  $\hat{\alpha}_i$ .

It remains to specify how to estimate the count at each time step from previously “released” noisy  $\mathbf{p}$ -sums. Let  $i_1 < i_2 < \dots < i_m$  denote the positions of non-zero bits in the binary representation of  $t$ . It is not hard to see that the count at time  $t$  is the sum of  $m$   $\mathbf{p}$ -sums of size  $2^{i_1}, 2^{i_2}, \dots, 2^{i_m}$  respectively. They correspond to the values of variables  $\alpha_{i_1}, \dots, \alpha_{i_m}$  maintained by the mechanism at time  $t$ . Specifically, the  $\mathbf{p}$ -sums are 1)  $\alpha_{i_1}$ , of the most recent



$2^{i_1}$  items; 2)  $\alpha_{i_2}$ , of the preceding  $2^{i_2}$  items; 3)  $\alpha_{i_3}$ , of the further preceding  $2^{i_3}$  items and so on.

**Implementing the Binary Mechanism with small memory.** Observe that the mechanism only needs to store  $\mathbf{p}$ -sums required for estimating the count at a future point in time, henceforth referred to as *active*  $\mathbf{p}$ -sums. The mechanism can safely discard  $\mathbf{p}$ -sums that are no longer needed, henceforth referred to as *inactive*  $\mathbf{p}$ -sums. For example, at  $t = 2^k$  where  $k \in \mathbb{N}$ , the only active  $\mathbf{p}$ -sum is  $\Sigma[1, t]$ . All other  $\mathbf{p}$ -sums computed between time 1 and  $t$  will no longer be needed after time  $t = 2^k$ . In the Binary Mechanism, we save the active  $\mathbf{p}$ -sums in variables  $\alpha_j$ 's. We reuse these variables, and let new and active  $\mathbf{p}$ -sums overwrite old and inactive ones. As a result, we only need to track  $O(\log t)$   $\mathbf{p}$ -sums and their noisy versions at time  $t$ .

To be more concrete, we describe what happens when a new item arrives at time  $t$  by making an analogy to a binary counter incremented from  $t - 1$  to  $t$ . Let  $i$  be the position of the least significant non-zero bit in  $t$ . Then  $t - 1$  has  $i - 1$  trailing 1's in its binary representation, and a carry occurs at position  $i$  when  $t - 1$  is incremented by 1. At the end of time  $t - 1$ , the mechanism stores  $\mathbf{p}$ -sums  $\alpha_j$  of sizes  $2^j$  for each  $j < i$ . During time step  $t$ , the mechanism performs the update  $\alpha_i \leftarrow \sum_{j < i} \alpha_j + \sigma(t)$ . Now  $\alpha_i$  is a new  $\mathbf{p}$ -sum of the most recent  $2^i$  items ending at time  $t$ . The mechanism also anonymizes  $\alpha_i$  with fresh randomness, and stores the corresponding noisy version  $\hat{\alpha}_i$ . Since the  $\mathbf{p}$ -sums  $\alpha_j$  for  $j < i$  are no longer needed, they (together with their noisy versions) are set to 0.

**Differential Privacy.** Consider an item arriving at  $t \in [T]$ . We analyze which of the  $\mathbf{p}$ -sums would be affected if  $\sigma(t)$  is flipped. It is not hard to see that the item  $\sigma(t)$  can be in at most  $\log T$   $\mathbf{p}$ -sums. In particular, it can be in at most 1  $\mathbf{p}$ -sum of size  $2^j$ , where  $j \leq \log T$ . Observe that each noisy  $\mathbf{p}$ -sum maintains  $\frac{\epsilon}{\log T}$ -differential privacy by Fact 1. Hence, we can conclude the differential privacy of the Binary Mechanism.

**Theorem 3** (Differential Privacy). *For  $T \in \mathbb{N}$ , the Binary Mechanism preserves  $T$ -bounded  $\epsilon$ -differential privacy.*

**Utility.** We next consider the usefulness of the Binary Mechanism. Each estimated count  $\mathcal{B}(t)$  is the sum of at most  $\log T$  noisy  $\mathbf{p}$ -sums, and each noisy  $\mathbf{p}$ -sum contains fresh, independent Laplace noise  $\text{Lap}(\frac{\log T}{\epsilon})$ . Therefore, the error at time  $t$  is the summation of at most  $O(\log t)$  i.i.d. Laplace distributions  $\text{Lap}(\frac{\log T}{\epsilon})$ . We use the Corollary 1 to conclude the mechanism's usefulness.

**Theorem 4** (Utility). *For each  $t \in [T]$ , the  $T$ -bounded Binary Mechanism is  $(O(\frac{1}{\epsilon}) \cdot (\log T) \cdot \sqrt{\log t} \cdot \log \frac{1}{\delta}, \delta)$ -useful at time  $t \in [T]$ .*

## 4 Unbounded Counting Mechanisms

Previously, we mainly considered time-bounded mechanisms, i.e., the mechanism requires a priori knowledge of an upper bound on the time. We now describe how to remove this assumption and derive unbounded counting mechanisms from the Binary Mechanism. The first approach, referred to as the modified Binary Mechanism, adds varying noise to each  $\mathbf{p}$ -sum.

The second approach, referred to as the Hybrid Mechanism, gives a generic way to convert any time-bounded mechanism  $\mathcal{M}$  into an unbounded one by running two mechanisms in parallel: (1) an unbounded mechanism that only outputs counts at time steps  $t$  being powers of 2; (2) a time-bounded mechanism  $\mathcal{M}$  to take care of items arriving at time steps between successive powers of 2.

### 4.1 Modified Binary Mechanism

The time-bounded Binary Mechanism needs to know  $T$ , an upper bound on time, in order to decide the magnitude of noise added to each  $\mathbf{p}$ -sum. How do we decide an appropriate magnitude of noise to add without knowledge of  $T$ ? One idea is to add varying noise to each  $\mathbf{p}$ -sum, depending on the length of the  $\mathbf{p}$ -sum.

**The modified Binary Mechanism.** The modified Binary Mechanism works the same way as the Binary Mechanism, except for the magnitude of noise we add to each p-sum. Suppose we add  $\text{Lap}(\frac{a_i}{\epsilon})$  noise to a p-sum consisting of  $2^i$  elements, where  $a_i > 0$ . In other words, replace Equation (3) in Figure 2 with the following:

$$\hat{\alpha}_i \leftarrow \alpha_i + \text{Lap}\left(\frac{a_i}{\epsilon}\right)$$

We leave the  $a_i$ 's as undetermined parameters for now, and work out the suitable values below.

**Theorem 5** (Differential Privacy). *For  $T \in \mathbb{N}$ , the modified Binary Mechanism preserves  $T$ -bounded  $(\epsilon \sum_{i \leq K} \frac{1}{a_i})$ -differential privacy, where  $K = \log_2 T + O(1)$ . Moreover, if the series  $\sum_i \frac{1}{a_i}$  converges to some value  $A < \infty$ , then the modified Binary Mechanism preserves  $\epsilon A$ -differential privacy.*

**Theorem 6** (Utility). *The modified Binary Mechanism is  $(O(\frac{1}{\epsilon} \cdot \sqrt{\sum_{i=0}^k a_i^2} \cdot \log \frac{1}{\delta}), \delta)$ -useful at time  $t$ , where  $k := \max\{i : \text{Bin}_i(t) = 1\} = \log_2 t + O(1)$ .*

**remark 1.** *Using Theorems 5 and 6, and scaling  $\epsilon$ , we can conclude that there is a  $T$ -bounded counting mechanism that preserves  $T$ -bounded  $\epsilon$ -differential privacy, and is  $(O(\frac{1}{\epsilon} \cdot \log \frac{1}{\delta}) \cdot \sum_{i=0}^K \frac{1}{a_i} \cdot \sqrt{\sum_{i=0}^K a_i^2}, \delta)$ -useful at time  $T$ , where  $K := \log_2 T + O(1)$ . Hence, loosely speaking, with fixed privacy parameter  $\epsilon$ , the error is  $\sum_{i=0}^K \frac{1}{a_i} \cdot \sqrt{\sum_{i=0}^K a_i^2}$ , and the best choice of  $a_i$ 's should minimize the above error term. In fact, the error term  $\sum_{i=0}^K \frac{1}{a_i} \cdot \sqrt{\sum_{i=0}^K a_i^2}$  is minimized when all  $a_i$ 's are equal, and this is exactly our time-bounded Binary Mechanism, where  $a_i = \log T$  for all  $1 \leq i \leq K$ .*

We next give different choices of  $a_i$ 's to obtain unbounded counting mechanisms with different guarantees.

**Corollary 3.** *Let  $0 < \delta < 1$ ,  $\epsilon > 0$ . Suppose  $\theta > 0$ , and set  $a_i := (i + 1)^{1+\theta}$ . There is an unbounded counting mechanism that preserves  $\epsilon$ -differential privacy and is  $(O(\frac{1}{\theta \epsilon} \cdot (\log t)^{1.5+\theta} \cdot \log \frac{1}{\delta}), \delta)$ -useful at time  $t$ .*

In Corollary 3, if we choose  $\theta > 0$  to be arbitrarily small, there will be a factor of  $\frac{1}{\theta}$  in the error. Instead, we choose  $a_i$  to be a function that is slightly super linear in  $i$ .

For  $n \in \mathbb{N}$ , define  $\kappa(n) := \max\{i : \ln^{(i)} n \geq 1\} = \Theta(\log^* n)$ . Recall that  $\log^*$  is defined in terms of base 2 logarithm.

For  $n \in \mathbb{N}$ , define the function  $\mathcal{L}(n) := \prod_{i=0}^{\kappa(n)} \ln^{(i)} n = n(\ln n)(\ln \ln n)(\ln \ln \ln n) \cdots$ .

**Claim 2.** *Let  $K \in \mathbb{N}$ . Then,  $\sum_{i=1}^K \frac{1}{\mathcal{L}(i)} \leq \Theta(\log^* K)$  and*

$$\sqrt{\sum_{i=1}^K \mathcal{L}(i)^2} \leq O(K^{1.5} \mathcal{L}(\log K)).$$

**Corollary 4.** *Let  $0 < \delta < 1$ ,  $\epsilon > 0$ . Set  $a_i := \mathcal{L}(i + 1)$ . There is an unbounded counting mechanism, such that for any  $t \in \mathbb{N}$ , it preserves  $O(\epsilon \log^*(\log t))$ -differential privacy, and is  $(O(\frac{1}{\epsilon}) \cdot \log^{1.5} t \cdot \mathcal{L}(\log \log t) \cdot \log \frac{1}{\delta}), \delta)$ -useful at time  $t$ .*

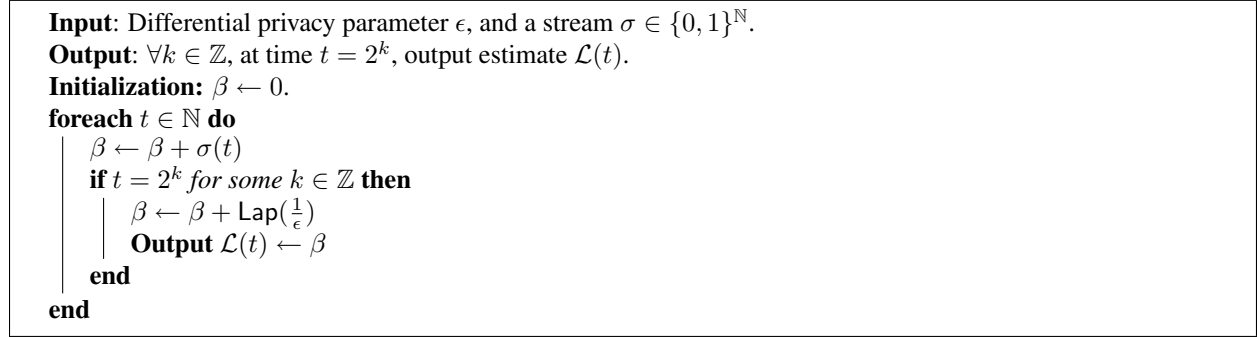
## 4.2 Hybrid Mechanism

The modified Binary Mechanism is unbounded, and achieves guarantees similar to the Binary Mechanism, but with a slight penalty in the error term (under fixed privacy parameter  $\epsilon$ ). We now describe a better approach that achieves time unboundedness and meanwhile provides the same asymptotic bounds as the Binary Mechanism.

The idea is to have an unbounded mechanism which only reports the estimated counts at sparse intervals, in particular, when  $t$  is a power of 2. We would expect such a mechanism to have better guarantees than one that has to report at every time step.

So what do we do when  $t$  is not a power of 2? We know the approximate count  $\hat{c}_1$  for the time period  $[1, T]$  where  $T = 2^k$  for some non-negative integer  $k$ . Suppose we also know the approximate count  $\hat{c}_2$  for the time period  $[T + 1, t]$  where  $T + 1 \leq t \leq 2T$ . Then we can estimate the count at time  $t$  as  $\hat{c}_1 + \hat{c}_2$ . Therefore, it remains for us to count the 1's between  $[T, t]$  for any  $t \in [T + 1, 2T]$ . We can simply apply a  $T$ -bounded mechanism (e.g., the Binary Mechanism) for this task.

**Logarithmic Counting Mechanism.** We now design an unbounded mechanism called the *Logarithmic Mechanism* which reports the count only when the time  $t$  is a power of 2.



**Figure 3:** Logarithmic Mechanism  $\mathcal{L}$

The idea for the Logarithmic Mechanism is quite simple. The mechanism internally keeps a value  $\beta$  which is initialized to 0.  $\beta$  is used to keep track of the approximate count at any point of time. As an item comes in, its value is added to  $\beta$ . At  $t$  equal to a power of 2, the mechanism adds fresh randomness to the value  $\beta$  (on top of randomness previously added), and outputs  $\beta$ .

If  $t$  is a power of 2, it is clear that the accumulated error at time  $t$  is a sum of  $O(\log t)$  independent Laplace distributions  $\text{Lap}(\frac{1}{\epsilon})$ . Hence, we have the following guarantee from Lemma 1.

**Theorem 7.** *The Logarithmic Counting Mechanism is unbounded, preserves  $\epsilon$ -differential privacy and is  $(O(\frac{1}{\epsilon}) \cdot \sqrt{\log t} \cdot \log \frac{1}{\delta}, \delta)$ -useful at time  $t = 2^k$  for some  $k \in \mathbb{Z}$ .*

**Logarithmic Mechanism: the p-sum view.** The Logarithmic Mechanism also has a p-sum interpretation. Equivalently, one can think of it as releasing the noisy p-sums  $\hat{\alpha} = \widehat{\Sigma}[1, 1]$ , as well as  $\hat{\beta}_k = \widehat{\Sigma}[2^k + 1, 2^{k+1}]$  for every  $k \in \mathbb{Z}$ . Now an observer can estimate the count at time  $t = 2^k$  as  $\hat{\alpha} + \sum_{i=0}^{k-1} \hat{\beta}_i$ .

**Hybrid Mechanism.** We combine the Logarithmic Mechanism and a time-bounded counting mechanism to process a given stream  $\sigma$ . We run one copy of  $\frac{\epsilon}{2}$ -differentially private Logarithmic Mechanism, which reports an approximate count when  $t$  is a power of 2. Suppose the Logarithmic Mechanism has reported count  $\mathcal{L}(T)$  at  $T = 2^k$  for some non-negative integer  $k$ . For time  $t$  in the range  $T + 1 \leq t \leq 2T$ , we run an  $\frac{\epsilon}{2}$ -differentially private  $T$ -bounded counting mechanism denoted as  $\mathcal{M}$  to count the number of 1's in the range  $[T + 1, t]$ . We write  $\tau = t - T$ . At time  $t$ , let  $\mathcal{M}(\tau)$  be the number of 1's in  $[T + 1, T + \tau]$  reported by the  $T$ -bounded counting mechanism  $\mathcal{M}$ . Then, the hybrid mechanism reports  $\mathcal{L}(T) + \mathcal{M}(\tau)$  at time  $t$ . The detailed Hybrid Mechanism is presented in Figure 4.

**Theorem 8.** *Assume that given any  $\epsilon > 0$  and  $0 < \delta < 1$ , Logarithmic Mechanism  $\mathcal{L}$  is  $\epsilon$ -differentially private and is  $(f(\epsilon, t, \delta), \delta)$ -useful at time  $t$ . Similarly, assume that given any  $\epsilon > 0$ , the  $T$ -bounded mechanism  $\mathcal{M}$  is  $\epsilon$ -differentially private and is  $(g(\epsilon, T, \tau, \delta), \frac{\delta}{2})$ -useful at time  $\tau \in [T]$ , where  $g$  is monotonically increasing with  $T$  and  $\tau$ . Then, the Hybrid Mechanism described above is unbounded, preserves  $\epsilon$ -differential privacy, and is  $(f(\frac{\epsilon}{2}, t, \frac{\delta}{2}) + g(\frac{\epsilon}{2}, t, t, \frac{\delta}{2}), \delta)$ -useful at time  $t$ .*

*Proof.* We first note that the  $\frac{\epsilon}{2}$ -differentially private Logarithmic Counting Mechanism is run in parallel with at most one instance of the  $\frac{\epsilon}{2}$ -differentially private time bounded counting mechanism at any time. Hence, the Hybrid Counting Mechanism is  $\epsilon$ -differentially private.

Suppose at time  $t = 2^k$  (for some non-negative integer  $k$ ), the Logarithmic Counting Mechanism has reported some count  $\mathcal{L}(t)$ . From the assumption, we know that with probability at least  $1 - \frac{\delta}{2}$ , the error of  $\mathcal{H}(t) := \mathcal{L}(t)$  is at most  $f(\frac{\epsilon}{2}, t, \frac{\delta}{2})$ .

Consider  $T + 1 \leq t \leq 2T$ . We write  $\tau := t - T$ . Suppose the  $T$ -bounded counting mechanism  $\mathcal{M}_T$  reports  $\mathcal{M}_T(\tau)$  at time  $t$ . From the assumption, with probability at least  $1 - \frac{\delta}{2}$ ,  $\mathcal{M}_T(\tau)$  has error at most  $g(\epsilon, T, \tau, \delta), \frac{\delta}{2}) \leq g(\frac{\epsilon}{2}, t, t, \frac{\delta}{2})$ .

We conclude that with probability at least  $1 - \delta$ ,  $\mathcal{H}(t) := \mathcal{L}(T) + \mathcal{M}_T(\tau)$  has error at most  $f(\frac{\epsilon}{2}, t, \frac{\delta}{2}) + g(\frac{\epsilon}{2}, t, t, \frac{\delta}{2})$ .  $\square$

```

Input: Differential privacy parameter  $\epsilon$ , a stream  $\sigma \in \{0, 1\}^{\mathbb{N}}$ , Logarithmic Mechanism  $\mathcal{L}$ , and a time-bounded mechanism  $\mathcal{M}$ .
Output: For each  $t \in \mathbb{N}$ , output estimate  $\mathcal{H}(t)$ .
Initialization:  $T \leftarrow 1$ .
Initiate the mechanism  $\mathcal{L}$  with privacy parameter  $\frac{\epsilon}{2}$  on stream  $\sigma$ .
foreach  $t \in \mathbb{N}$  do
  Feed  $\sigma(t)$  to mechanism  $\mathcal{L}$ .
  if  $t = 2^k$  for some  $k \in \mathbb{Z}$  then
    Output  $\mathcal{H}(t) \leftarrow \mathcal{L}(t)$ 
     $T \leftarrow t$ 
    Initiate an instance of the  $T$ -bounded mechanism  $\mathcal{M}_T$  with time upper bound  $T$ , privacy parameter  $\frac{\epsilon}{2}$  and stream  $\sigma^{(T)} \in \{0, 1\}^T$ , where  $\sigma^{(T)}(\tau) := \sigma(\tau + T)$  for  $\tau \in [1, T]$ .
  else
     $\tau \leftarrow t - T$ 
    Feed  $\sigma^{(T)}(\tau) := \sigma(t)$  to mechanism  $\mathcal{M}_T$ .
    Output  $\mathcal{H}(t) \leftarrow \mathcal{L}(T) + \mathcal{M}_T(\tau)$ 
  end
end

// At time  $t$ ,  $T$  is the largest power of 2 no bigger than  $t$ .
//  $\sigma^{(T)}$  is the sub-stream of  $\sigma$  for the duration  $[T+1, 2T]$ .
//  $\mathcal{M}_T$  is a time-bounded mechanism that runs for  $[T+1, 2T]$ .

```

**Figure 4:** Hybrid Mechanism  $\mathcal{H}$  (with Mechanism  $\mathcal{M}$ ).

**Corollary 5.** *If we instantiate the Hybrid Mechanism using the Binary Mechanism as the  $T$ -bounded mechanism, the resulting Hybrid Mechanism is unbounded, preserves  $\epsilon$ -differential privacy, and is  $(O(\frac{1}{\epsilon}) \cdot (\log t)^{1.5} \cdot \log \frac{1}{\delta}, \delta)$ -useful at time  $t$ .*

For simplicity, in the remainder of the paper, when we refer to the Hybrid Mechanism, we mean the Hybrid Mechanism instantiated with the Binary Mechanism.

**Hybrid Mechanism: the p-sum view.** One can also interpret the Hybrid Mechanism naturally using the p-sum framework. Basically, one can equivalently think of the Hybrid Mechanism as releasing the union of the noisy p-sums of the Logarithmic Mechanism and the Binary Mechanism. From this set of noisy p-sums, an observer can compute the approximate count at every time step  $t \in \mathbb{N}$ .

## 5 Achieving Pan Privacy

The mechanisms described thus far are not designed to resist intrusions in which the adversary can learn snapshots of the mechanism's internal states. We now consider how to add pan privacy, i.e., the ability to resist intrusions into the mechanisms.

### 5.1 Pan Privacy Against Single Unannounced Intrusion

Recall that all our mechanisms fall within the p-sum framework described in Section 3.2. We now describe a generic technique for transforming any  $\epsilon$ -differentially private mechanism in the p-sum framework into an  $\epsilon$ -pan private mechanism resilient against a single unannounced intrusion, with only a constant-factor loss in the mechanism's utility. The techniques used here are similar to those proposed by Dwork *et.al.* [7, 10].

A mechanism in the p-sum framework releases a sequence of noisy p-sums from which an observer can estimate the count at each time step. One way to implement such a mechanism is the following:

```

Let  $P$  denote the (possibly infinite) set of all noisy p-sums the mechanism intends to output.
Initialize:  $S \leftarrow \emptyset$ .
//  $S$  is the set of currently active p-sums
foreach  $t \in \mathbb{N}$  do
  // Initialize a counter for each
  // relevant p-sum starting at  $t$ :
  foreach noisy p-sum  $\widehat{\Sigma}[t, t'] \in P$  do
     $\beta_{t, t'} \leftarrow 0$ 
     $S \leftarrow S \cup \{\beta_{t, t'}\}$ 
  end
  foreach  $\beta_{t_1, t_2} \in S$  do
    // Add the current item to all relevant counters
    if  $t_1 \leq t \leq t_2$  then
      |  $\beta_{t_1, t_2} \leftarrow \beta_{t_1, t_2} + \sigma(t)$ 
    end
    // Output noisy p-sums ending at  $t$ 
    // and remove these counters from memory
    if  $t = t_2$  then
      | Output  $\widehat{\Sigma}[t_1, t_2] = \beta_{t_1, t_2} + \text{Lap}(\frac{a}{\epsilon})$  with appropriate choice of  $a$ 
      |  $S \leftarrow S \setminus \{\beta_{t_1, t_2}\}$ 
    end
  end
end

```

(5)

Basically, if  $\widehat{\Sigma}[t_1, t_2]$  is a noisy p-sum the mechanism intends to release, then the mechanism initializes a corresponding counter  $\beta_{t_1, t_2}$  at time  $t_1$ . For  $t_1 \leq t \leq t_2$ ,  $\beta_{t_1, t_2} = \Sigma[t_1, t]$  at the end of  $t$ . In other words, the mechanism internally keeps track of the accurate count for the duration  $[t_1, t]$ . At the end of time  $t_2$ ,  $\beta_{t_1, t_2} = \Sigma[t_1, t_2]$ . Now the mechanism adds noise to the counter, and reveals the noisy p-sum  $\widehat{\Sigma}[t_1, t_2]$ .

**remark 2.** *If we implement the Binary Mechanism or the Hybrid Mechanism using the above approach, it is not hard to show that only  $O(\log t)$  counters are needed at time  $t$ . In other words, at any time  $t$ ,  $|S| < O(\log t)$ .*

We now examine how our current mechanisms fail to achieve pan privacy, and show how to remedy the problem. Notice that if an intrusion happens at time  $t_1 \leq t < t_2$ , then the adversary can learn  $\Sigma[t_1, t]$  and privacy is obviously broken. The way to resolve this is to initialize the counter with some noise  $\text{Lap}(\frac{1}{\epsilon})$ . Basically, replace Equation (5) with the following for an appropriate choice of  $a$ .

$$\beta_{t, t'} \leftarrow \text{Lap}\left(\frac{a}{\epsilon}\right)$$

Specifically, we choose the magnitude of noise  $a$  as below. Suppose the original mechanism intended to add  $\text{Lap}(\frac{a(t_1, t_2)}{\epsilon})$  to the p-sum  $\Sigma[t_1, t_2]$ . Then we initialize the counter  $\beta_{t_1, t_2}$  with noise  $\text{Lap}(\frac{a(t_1, t_2)}{\epsilon})$ .

As a result of the above modification, each noisy p-sum output by the mechanism now has two independent Laplace noises, one added at the time the counter was initialized, the other added at the time of output. However, this would only lead to a constant factor increase in the error bound.

Intuitively, the above modification achieves  $\epsilon$ -pan privacy due to the following observations. Suppose the intrusion happens at time  $t_1 \leq t \leq t_2$ , the adversary learns the internal state  $\beta_{t_1, t_2} = \Sigma[t_1, t] + \text{Lap}(\frac{a}{\epsilon})$  in addition to the noisy p-sum  $\widehat{\Sigma}[t_1, t_2] = \Sigma[t_1, t_2] + \text{Lap}(\frac{a}{\epsilon}) + \text{Lap}(\frac{a}{\epsilon})$  output at time  $t_2$ . This is equivalent to revealing the noisy p-sums  $\widehat{\Sigma}[t_1, t] := \Sigma[t_1, t] + \text{Lap}(\frac{a}{\epsilon})$  and  $\widehat{\Sigma}[t+1, t_2] := \Sigma[t+1, t_2] + \text{Lap}(\frac{a}{\epsilon})$  to the adversary. Note that changing any single

position in the stream can only affect at most one of  $\widehat{\Sigma}[t_1, t]$  and  $\widehat{\Sigma}[t + 1, t_2]$ . It is not hard to see that this mechanism achieves  $\epsilon$ -pan privacy if its non-pan private version achieves  $\epsilon$  differential privacy.

Hence, we can conclude that all our counting mechanisms which fall within the p-sum framework can be converted to pan private versions with the same asymptotic guarantees.

## 5.2 Pan Privacy Against Multiple (Afterwards) Announced Intrusions

We can apply the same idea for multiple intrusions, as long as the mechanism is aware of the intrusions immediately after they are made. We initialize each counter  $\beta_{t_1, t_2}$  with fresh Laplace noise  $\text{Lap}(\frac{a}{\epsilon})$ . Whenever we detect an intrusion, we add fresh Laplace noise  $\text{Lap}(\frac{a}{\epsilon})$  to each active counter in memory. Finally, we add yet another fresh Laplace noise  $\text{Lap}(\frac{a}{\epsilon})$  before a p-sum is output by the mechanism.

One can see that the mechanism preserves pan privacy with the same guarantee as before. Moreover, the error scales with a square root factor on the number of intrusions made. In particular, we have the following theorem.

**Theorem 9.** *If we modify the Hybrid Mechanism using the above approach, the resulting mechanism is unbounded, preserves  $\epsilon$ -pan privacy against multiple intrusions, and is  $(O(\frac{\sqrt{k}}{\epsilon}) \cdot (\log t)^{1.5} \cdot \log \frac{1}{\delta}, \delta)$ -useful at time  $t$ , where  $k$  is the number of intrusions made before time  $t$ .*

## 6 Applications

The counter is a fundamental primitive in many streaming algorithms [2, 15, 17]. We believe that our differentially private continual counter construction will inspire the design of a wide class of differentially private streaming algorithms.

Our counting mechanism also has immediate practical applications. For example, it is a common practice for websites such as IMDB, Delicious and Digg to suggest the most popular movies, bookmarks or news items to visitors. Our construction can be extended to suggest top- $k$  items or hot items in a differentially private manner.

**Continual top- $k$  suggestions.** Let  $\mathcal{U} := \{1, 2, \dots, m\}$  denote the universe of possible items. Let  $\sigma \in \mathcal{U}^{\mathbb{N}}$  denote a stream of incoming items. Specifically,  $\sigma(t) \in \mathcal{U}$  denotes the item at time  $t \in \mathbb{N}$ . At every time step  $t \in \mathbb{N}$ , the continual top- $k$  mechanism  $\mathcal{K}$  outputs a tuple  $\mathcal{K}(\sigma)(t) \in \mathcal{U}^k$ , an approximation of the top- $k$  items seen thus far.

We say that two streams  $\sigma, \sigma'$  are adjacent if they differ at exactly one time step  $t \in \mathbb{N}$ . A continual top- $k$  mechanism is  $\epsilon$ -differentially private if for any adjacent streams  $\sigma$  and  $\sigma'$ , and any measurable subset  $S \subseteq (\mathcal{U}^k)^{\mathbb{N}}$ ,  $\Pr[\mathcal{K}(\sigma) \in S] \leq \exp(\epsilon) \cdot \Pr[\mathcal{K}(\sigma') \in S]$ .

We can use the continual counter construction to give continual top- $k$  suggestions in a differentially private manner. To achieve this, the mechanism  $\mathcal{K}$  internally runs  $m$  differentially private continual counters, to keep track of the approximate count of each item thus far. We assume that each continual counter is instantiated using the Hybrid Mechanism described in Section 4.2. At each time step, the mechanism  $\mathcal{K}$  ranks the items according to their approximate counts as suggested by the  $m$  counters, and outputs the resulting top- $k$ .

More formally, we “split” the stream  $\sigma$  into  $m$  different indicator streams  $\sigma_1, \sigma_2, \dots, \sigma_m$ . If  $\sigma(t) = i \in \mathcal{U}$ , then  $\sigma_i(t) := 1$  and  $\sigma_j(t) := 0$  for all  $j \neq i$ . The mechanism  $\mathcal{K}$  instantiates  $m$   $\epsilon$ -differentially private continual counters  $\mathcal{H}_1, \dots, \mathcal{H}_m$  using the Hybrid Mechanism, and feeds the indicator stream  $\sigma_i$  to the counter  $\mathcal{H}_i$  for each  $i \in \mathcal{U}$ . At time  $t$ , the mechanism  $\mathcal{K}$  computes the ranking according to the approximate counts  $\{\mathcal{H}_1(t), \mathcal{H}_2(t), \dots, \mathcal{H}_m(t)\}$ , and outputs the resulting top- $k$  items.

The mechanism  $\mathcal{K}$  described above preserves  $2\epsilon$ -differential privacy, and achieves good accuracy in the following sense.

**Claim 3.** *The above continual top- $k$  mechanism  $\mathcal{K}$  is  $2\epsilon$ -differentially private. Furthermore, at any time  $t \in \mathbb{N}$ , for any two items  $i, j \in \mathcal{U}$  if the true counts of items  $i$  and  $j$  differ by  $\Omega(\frac{1}{\epsilon} \cdot (\log t)^{1.5} \cdot \log \frac{1}{\delta})$ , then with probability  $1 - \delta$ , the ranking computed by mechanism  $\mathcal{K}$  preserves the correct ordering of  $i$  and  $j$  at time  $t$ .*

*Proof. (sketch.)* For  $2\epsilon$  differential privacy, observe that flipping one item in the original stream  $\sigma$  affects two of the counters among  $\mathcal{H}_1, \dots, \mathcal{H}_m$ . The usefulness argument follows directly from Corollary 5.  $\square$

**Continual hot items suggestions.** Hot items can be considered as a variation of top- $k$  queries. Using a similar idea, we can design a mechanism that continually outputs hot items suggestions in a differentially private manner. Basically, we run  $m$  differentially private continual counters to keep track of the approximate count for each item, and compute the hot items list according to their approximate counts. In this way, we can achieve the same guarantees as suggested by Claim 3.

Finally, note that the mechanisms described in this section can be augmented to achieve pan privacy using techniques described in Section 5.

## 7 Conclusion and Open Problems

We consider how a website or data aggregator can privately and continually release new statistics when the database is evolving over time. We propose an  $\epsilon$ -differentially private continual counter which has only poly-log error with respect to the time. One problem left open by this study is how to lower bound the error for the  $\epsilon$ -differentially private continual counter.

This represents an exciting and important new setting for differential privacy, as numerous websites adopt the practice of releasing new user statistics over time. A promising direction for future work is the design of a broad class of private streaming algorithms which continually output data. Another interesting direction is to implement the theoretic results and study their feasibility in real-world applications.

## Acknowledgments

We would like to thank John Bethencourt, Arvind Narayanan, Aaron Roth and Emil Stefanov for helpful discussions and insightful feedback.

## References

- [1] J. A. Calandrino, A. Narayanan, E. W. Felten, and V. Shmatikov. Don't review that book: Privacy risks of collaborative filtering. Manuscript, 2009.
- [2] E. D. Demaine, A. López-Ortiz, and J. I. Munro. Frequency estimation of internet packet streams with limited space. In *ESA '02: Proceedings of the 10th Annual European Symposium on Algorithms*, 2002.
- [3] I. Dinur and K. Nissim. Revealing information while preserving privacy. In *PODS*, 2003.
- [4] C. Dwork. Differential privacy. Invited talk at *ICALP*, 2006.
- [5] C. Dwork. Differential privacy: A survey of results. In *TAMC*, 2008.
- [6] C. Dwork. The differential privacy frontier. In *TCC*, 2009.
- [7] C. Dwork. Differential privacy in new settings. Invited presentation at *ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 2010.
- [8] C. Dwork. A firm foundation for private data analysis. In *Communications of the ACM*, 2010.
- [9] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In *TCC*, 2006.
- [10] C. Dwork, M. Naor, T. Pitassi, and G. N. Rothblum. Differential privacy under continual observation. Manuscript in preparation, 2009.
- [11] C. Dwork, M. Naor, T. Pitassi, G. N. Rothblum, and S. Yekhanin. Pan-private streaming algorithms. In *Innovations in Computer Science (ISC)*, 2010.

- [12] C. Dwork and S. Yekhanin. New efficient attacks on statistical disclosure control mechanisms. In *CRYPTO*, 2008.
- [13] R. Jones, R. Kumar, B. Pang, and A. Tomkins. Vanity fair: privacy in querylog bundles. In *CIKM*, 2008.
- [14] A. Korolova, K. Kenthapadi, N. Mishra, and A. Ntoulas. Releasing search queries and clicks privately. In *WWW*, 2009.
- [15] G. S. Manku and R. Motwani. Approximate frequency counts over data streams. In *VLDB*, 2002.
- [16] F. McSherry and I. Mironov. Differentially private recommender systems: building privacy into the netflix prize contenders. In *KDD*, 2009.
- [17] A. Metwally, D. Agrawal, and A. E. Abbadi. Efficient computation of frequent and top-k elements in data streams. In *ICDT*, 2005.
- [18] A. Narayanan and V. Shmatikov. Robust de-anonymization of large sparse datasets. In *IEEE Symposium on Security and Privacy*, 2008.
- [19] S. L. Warner. Randomized response: A survey technique for eliminating evasive answer bias. *Journal of the American Statistical Association*, 1965.