

How to Make Private Distributed Cardinality Estimation Practical, and Get Differential Privacy for Free[†]

Changhui Hu¹, Jin Li², Zheli Liu^{3,‡}, Xiaojie Guo³, Yu Wei³, Xuan Guang⁴, Grigorios Loukides⁵
Changyu Dong^{1,‡}

¹ *School of Computing, Newcastle University, {changhui.hu, changyu.dong}@newcastle.ac.uk*

² *Institute of AI and Blockchain, Guangzhou University, lijn@gzhu.edu.cn*

³ *College of Cyber Science, Nankai University, liuzheli@nankai.edu.cn
{xiaojie.guo, stoneboat}@mail.nankai.edu.cn*

⁴ *School of Mathematical Sciences and LPMC, Nankai University, xguang@nankai.edu.cn*

⁵ *Department of Informatics, King's College London, grigorios.loukides@kcl.ac.uk*

Abstract

Secure computation is a promising privacy enhancing technology, but it is often not scalable enough for data intensive applications. On the other hand, the use of sketches has gained popularity in data mining, because sketches often give rise to highly efficient and scalable sub-linear algorithms. It is natural to ask: what if we put secure computation and sketches together? We investigated the question and the findings are interesting: we can get security, we can get scalability, and somewhat unexpectedly, we can also get differential privacy – for free. Our study started from building a secure computation protocol based on the Flajolet-Martin (FM) sketches, for solving the Private Distributed Cardinality Estimation (PDCE) problem, which is a fundamental problem with applications ranging from crowd tracking to network monitoring. The state of art protocol for PDCE (Fenske et al. CCS'17) is computationally expensive and not scalable enough to cope with big data applications, which prompted us to design a better protocol. Our further analysis revealed that if the cardinality to be estimated is large enough, our protocol can achieve (ϵ, δ) -differential privacy automatically, without requiring any additional manipulation of the output. The result signifies a new approach for achieving differential privacy that departs from the mainstream approach (i.e. adding noise to the result). Free differential privacy can be achieved because of two reasons: secure computation minimizes information leakage, and the intrinsic estimation variance of the FM sketch makes the output of our protocol uncertain. We further show that the result is not just theoretical: the minimal cardinality for differential privacy to hold is only $10^2 - 10^4$ for typical parameters.

[†]The full version of this paper can be found here: <https://eprint.iacr.org/2020/1576>

[‡]Changyu Dong and Zheli Liu are the corresponding authors.

1 Introduction

Data privacy has become an increasingly acute problem, especially when the hunger for data drives large-scale collection and (mis)use, without well-thought-out precautions in place. The tension between data utilization and data privacy has developed into a societal challenge and led to stricter regulations, such as HIPAA [1], GLBA [2] and GDPR [3]. The pressing need for privacy has greatly stimulated research on secure computation [16]. Secure computation allows collaborative computation over private datasets held by multiple mutually untrusted parties, without revealing any information except what can be inferred from the output. Thus, secure computation has been regarded as one of the key privacy enhancing technologies [4].

While many secure computation protocols have been proposed to carry out various data processing tasks in a privacy preserving fashion, their scalability is often open to doubt. Despite the fact that the efficiency of secure computation has been drastically improved, secure computation is still orders of magnitude slower than computation in the clear. The overhead might be acceptable if the data to be processed is small, but it can be prohibitive when the data is big. Yet, the “killer” applications of secure computation are often data-intensive, and this has become a major impediment to the widespread use of secure computation.

One good example is *Private Distributed Cardinality Estimation* (PDCE). Cardinality estimation, the task of determining the number of distinct elements in the union of multiple sets, is of particular importance in databases, data mining and distributed systems [5, 34, 35, 62]. While the task is easy to perform when data is in a single small database, it becomes challenging when data is collected independently from multiple sources at a high rate [48]. Naively maintaining a counter

at each source and summing the counters up will not work because more often than not, there are duplicates in the data being collected. The task is even more challenging if privacy is needed. PDCE has numerous applications, for example:

- **Scientific research and user studies.** Surveys and questionnaires are commonly used in medical science, social science and business studies to help researchers discover interesting correlations (e.g. [52]). It is not uncommon that several organizations independently collect data through surveys and questionnaires over the same population. For example, diet habits could be surveyed by researchers from a medical institution, a government agency, and an insurance company, for different projects. If pooled together, the data would be of a much higher utility and could lead to improved decision making. For instance, the number of distinct individuals across all datasets having a certain diet habit could help identifying risk factors related to chronic diseases such as diabetes, while each individual dataset may be too small to draw a convincing conclusion. However, the data cannot be pooled together in practice because of privacy obligations imposed on each data collector.
- **Crowd counting and tracking** [61]. The task of estimating the number of individuals entering or passing by a place is key in urban planning, surveillance, public health study and retail analytics, to understand the effectiveness of building and road design, the patterns of human mobility, the spread of infectious diseases, and the patterns of customer behavior. Many existing commercial systems identify and track people through personally identifiable information (PII), such as fingerprints of mobile devices [9] or MAC addresses of WiFi cards [47, 49], collected through a distributed network of sensors or WiFi hot-spots that are deployed e.g. in a big shopping mall or a retail chain across the country [57]. Usually, to obtain the estimate, the data is transmitted to, stored and processed in a central database. However, this has already raised widespread privacy concerns [58, 59]. Ideally the estimate should be obtained without the need to store or transmit PII. Also, if published, it should not disclose information about any specific individual.
- **Network monitoring and statistics.** An example is the detection of DDoS attacks by collecting information at an ISP's border routers and identifying sudden increases in the total number of distinct source IP addresses. This detection method works because the attacker usually commands many "zombies" distributed across the Internet to send packets with randomly spoofed IP source addresses to the victim [46]. Another example is that many websites nowadays use Content Distribution Networks to provide load balancing and fast access. One statistic that web masters often want to estimate is the total number, across all replicas, of distinct visitors who accessed the website [7]. This can be reduced to the distributed cardinality estimation

problem. In both examples, privacy concerns are raised due to the fact that IP addresses can be used to track back users, revealing confidential information ranging from their location to personal and behavioral traits. A second concern, posed by the scale of the Internet, is to be able to address the distributed cardinality estimation problem efficiently.

Driven by the need, PDCE has attracted substantial interest [8, 17, 20, 23, 27, 28, 30, 38, 43, 56, 60]. The current state of the art is a secure computation protocol proposed by Fenske et.al. in CCS'17 [30]. Functionality and privacy-wise, the protocol is impeccable. However, it is not scalable enough, because it relies on expensive public key encryption and computationally demanding sub-protocols such as verifiable shuffling. The running time of the protocol is in the order of hours when the cardinality is in the order of 10^4 . The protocol can fulfill its task for measuring the number of visitors to the Tor network because the cardinality to be estimated is small. However, it cannot cope with mainstream big data applications, involving million or billion sized sets, because the protocol would require weeks or years to finish.

In the data mining community, the use of sketches has gained popularity recently [15]. Sketches are space-efficient data structures that summarize massive data, so that it can be efficiently processed, stored, and queried. Sketches allow representing data in sub-linear or constant space, and thus can be employed to improve the efficiency and scalability of algorithms. Sketches are lossy and do not preserve all the information in the data they represent. Thus, sketch-based computation returns only approximate answers. That said, big data applications often do not require exact answers and the parameters of sketches can often be adjusted to obtain sufficiently accurate answers. Due to the use of sketches, many real world systems can keep up with exponentially increasing data (e.g. [6, 37]).

Contributions In this paper, we build and analyse a new secure computation protocol based on the Flajolet-Martin (FM) sketch [31], for solving the PDCE problem. Initially, our intention was to make PDCE more practical. The protocol fulfils this intention very well. As expected, the protocol achieves high security, as well as much better efficiency and scalability than the state of the art [30]. Yet, this is not the end of the story. In the study we also found that the combination of secure computation and the FM sketch allows us to obtain differential privacy at no extra cost. This is interesting because differential privacy is a much desired property in PDCE, and neither secure computation nor the FM sketch provides it on its own. In more detail, the protocol has the following important features:

- **Highly Secure** Similar to existing work [27, 28, 30, 38, 43], we consider the scenario where a set of Data Parties (DPs) collect data and want to use some *untrusted* Computation Parties (CPs) to aggregate the data and estimate cardinality.

In such a setting, the untrusted CPs are modelled as corrupted and controlled by a single adversary. Unlike the vast majority of previous work that considers only semi-honest CPs, our protocol is developed on top of the SPDZ framework [19], and thus is secure in the presence of malicious CPs that can behave arbitrarily. For the total of d CPs, our protocol can tolerate up to $d - 1$ corrupted malicious CPs. Our protocol can achieve the following security goals as long as there exists one honest CP: (1) the adversary learns nothing from executing the protocol except the output of the protocol; (2) the adversary cannot affect the correctness of the computation without being detected. We formally prove the security of the protocol in the UC model [12].

- **Efficient and scalable.** We design our protocol around the FM sketch. By using FM sketches, we can accurately estimate the cardinality, while reducing the complexity of our protocol to logarithmic (in the maximum cardinality to be estimated). This is in contrast to the majority of the existing protocols [8, 17, 20, 27, 28, 30, 38, 43, 56] whose complexity is linear. Also, we reduce much of the computation needed for the online phase by using offline pre-processing. As a comparison, the protocol in [30] needs almost 1 hour in LAN to estimate the cardinality of a set containing **30,000** elements; while our protocol only needs less than 50 minutes (in WAN) to estimate the cardinality of a set containing **1 billion** elements, and the online phase running time is only about 5 seconds.
- **Offers differential privacy for free.** The most interesting finding from our study is that our protocol can achieve differential privacy [25] for free (i.e. without the need to add noise and/or further manipulate the output). In the security models for secure computation, the adversary is allowed to infer information from the output of the protocol. This sometimes is inadequate because individuals may still be re-identified through such inference. It is often desirable in applications like PDCE to additionally disallow such inference attacks by making the output from the protocol differentially private. We proved that, given the privacy parameters (ϵ, δ) and FM sketch parameters, if the cardinality to be estimated is sufficiently large, then the estimated cardinality output from our protocol satisfies (ϵ, δ) -differential privacy. What makes the finding so interesting is that neither secure computation nor sketches provide differential privacy on their own. However, we showed for the first time that when we put the two together, they complement each other by providing something the other lacks. The intrinsic estimation variance of FM sketches now makes the output of secure computation uncertain, thus can substitute the noise we usually need to add in order to achieve differential privacy. Secure computation makes it possible to do the computation without revealing anything except the output, which means the sketches are now hidden and any information leaked by the sketches is now concealed.

As a consequence, differential privacy can be achieved. We further show that this is not just a theoretical result. The lower bound of the cardinality for differential privacy to hold is reasonably small. Given typical parameters, the lower bound is usually only $10^2 - 10^4$. Thus, differential privacy can be easily satisfied in real world applications with our protocol. The technique used in our analysis is quite general, thus we would expect that with some modifications, it could be applied to other sketch based secure computation protocols as well. As the last remark, existing PDCE protocols [28, 30, 38, 43, 56], achieve differential privacy by adding noise, which however incurs a cost. This is especially true in [30], in which a large portion of the computation is spent on encrypting a large number ($10^4 - 10^5$) of noise bits and shuffling them with the data. Therefore, free differential privacy is beneficial to the efficiency and scalability of our protocol as well.

2 Related Work

In the literature, several PDCE protocols are also called Private Set Union Cardinality (PSU-CA) protocols [17, 20, 23, 27, 30]. However, the original definition of PSU-CA [17] requires the output to be the *exact* cardinality, while quite a few protocols [23, 27, 30] output an *estimate* close to the exact cardinality. To avoid confusion, we use the term PDCE in this paper and regard PSU-CA as a special case (in which the estimation error is 0). Note that not outputting the exact cardinality is not necessarily a deficiency. When differential privacy is required, the output anyway cannot be exact.

There are two different flavours of PDCE protocols: the first is that the DPs collect and compute the cardinality, without using CPs; the second is that the DPs only collect the data, and the CPs compute the estimation. We call the former DP-PDCE and the latter CP-PDCE to differentiate them. Our protocol is a CP-PDCE protocol. One approach [17, 20] for DP-PDCE is to reduce it to a Private Set Intersection Cardinality (PSI-CA) problem. The cardinality of union can be obtained by using the inclusion-exclusion principle. However, the inclusion-exclusion principle leads to exponential complexity (in the number of sets), therefore those protocols are limited to the two-party case. There are a few DP-PDCE [20] and CP-PDCE [8, 27] protocols based on Bloom Filters. The protocol in [8] is not secure, and [27] proposed a more secure variant of the protocol. The protocol in [20], as mentioned earlier, uses the inclusion-exclusion principle, thus is not scalable. All the above protocols have computational and communication complexity linear in the maximum cardinality to be estimated. FM sketches were used by the DP-PDCE protocol in [23] to lower its complexity to logarithmic. However, only a two-party protocol was given in the paper with a brief statement that a multiparty protocol is feasible. The DP-PDCE protocol in [60] also uses FM sketches. However this protocol is not secure. The protocol reveals more information

than the cardinality itself because the parties learn the union sketch in the protocol. It also assumes none of the parties collude, which is a very strong assumption. None of the aforementioned protocols supports differential privacy. There are protocols that provide differential privacy [28, 30, 38, 43, 56]. The CP-PDCE protocols in [28, 30, 38, 43] were all designed for gathering statistics in the Tor network [22], which naturally requires a high degree of privacy as the aim of Tor is to keep users anonymous. The protocols in [28, 38, 43] consolidate the observations of each DP into a counter, thus cannot eliminate duplicates when the counters are aggregated together. In [30], each DP maintains a hash-table with a public hash function for the observations. If an observation occurs multiple times, regardless by the same DP or by different DPs, it will be hashed into the same bin of the hash-table and the duplicates can be eliminated. However, in order to reduce collisions and maintain a reasonable accuracy, the hash-table size needs to be much larger than the maximum cardinality to be estimated. This impacts the efficiency and scalability of the protocol significantly. In [56], each DP represents its observations as a bit vector, enforces differential privacy on the vector using randomized response, and then passes the vector to a CP who can estimate cardinality of the set union. The estimation has a high standard deviation (in the order of the size of the universe of the set), thus the result is not accurate enough for many applications. All the above protocols except [30] consider the semi-honest or an even weaker adversary model, mainly for efficiency reasons, while our protocol and [30] are secure against more powerful malicious adversaries.

There is a large body of research works on Private Data Aggregation in which multiple data collectors (DPs) and data aggregators (CPs) are involved in aggregating data and outputting some statistics. Some works consider a much weaker security model and assume a trusted aggregator, who aggregates data from the DPs in plaintext and then adds noise before outputting the result [45, 55]. There are protocols that consider an untrusted aggregator, e.g. for computing private sum [14, 51, 54], or for frequency estimation over categorical data [14, 29], or for computing KNN and median [44]. Sketches (e.g. Count and Count-min sketches) were used in [44, 45] to make the protocols more efficient.

3 Preliminaries

3.1 Flajolet-Martin (FM) Sketches

We briefly review FM sketches. More details and analysis can be found in [23, 31, 53]. An FM sketch is a probabilistic data structure for counting the number of distinct elements in a multi-set. The data structure is a w -bit binary vector. Let FS denote an FM sketch, and FS[i] ($0 \leq i \leq w-1$) denote the i th bit in FS. An FM sketch is built using two functions:

- $H : \{0, 1\}^* \rightarrow \{0, 1\}^{w-1}$: a hash function that maps an input uniformly to a $(w-1)$ -bit string.
- $\rho : \{0, 1\}^{w-1} \rightarrow [0, w-1]$: a function that takes a $(w-1)$ -bit string as input and returns the number of trailing zeroes in it.

Initially, all bits in FS are set to 0. To estimate the cardinality of a multi-set S , for each element $x \in S$, we hash x and set FS[$\rho(H(x))$] = 1. The quantity N , which is the number of distinct elements in S , can be estimated using an estimator z_N that is the index of the first¹ 0 bit in FS, i.e. FS[z_N] = 0 and $\forall 0 \leq j < z_N$, FS[j] = 1. The expected value of z_N is close to $\log(\phi N)$, where $\phi = 0.77351$ is a correction factor. Therefore, N is roughly $2^{z_N}/\phi$. It is clear that the size of the sketch w must be larger than $\log(\phi N)$, otherwise z_N might not be correct. As suggested in [31], $w \geq \log(N) + 4$ should suffice.

The standard deviation of z_N is 1.12, which is too high (i.e. an estimation using z_N will typically be one binary order of magnitude off the true cardinality). To remedy this problem, [31] suggested to use m sketches, each with an independent hash function. Then we can obtain m estimators $z_{N,1}, \dots, z_{N,m}$, sum them to $Z_N = z_{N,1} + \dots + z_{N,m}$, and use the average $\frac{Z_N}{m}$ to estimate the cardinality N . The standard deviation of Z_N is $1.12 \cdot \sqrt{m}$. Thus, the standard deviation of $\frac{Z_N}{m}$ is $\frac{1.12}{\sqrt{m}}$, which is much smaller. In [53], the authors suggested the following, modified formula that can achieve better estimation accuracy:

$$\tilde{N} = \frac{2^{\frac{Z_N}{m}} - 2^{-\kappa \frac{Z_N}{m}}}{\phi} \quad (1)$$

where \tilde{N} is the cardinality estimated from m sketches, and $\kappa = 1.75$ is a correcting factor. In [23], it was shown that the accuracy of the estimation can be improved by enlarging m . This implies that the accuracy of the estimation can be adjusted to the desired level, by choosing a suitable m .

An important property of FM sketches that we use in the design of our protocol is that they can be merged. If we have two FM sketches FS₁ and FS₂ built with the same hash function, but on different sets S_1 and S_2 respectively, then bit-wisely ORing the two sketches produces a new FM sketch FS_U that counts the union of the two sets S_1 and S_2 . This process is lossless: FS_U is exactly the same as the sketch built using the union from the scratch. This holds also in the case of more than two sketches. Our protocol will use this property to union FM sketches from different DPs.

3.2 SPDZ

In this section, we briefly review the SPDZ scheme [18, 19, 40, 41] that will be used as the underlying framework for our protocol. We will follow mostly the notations in [40, 41]. Essentially, SPDZ is a secret-sharing based multiparty computation (MPC) scheme that supports secure computation over a

¹We use the most significant bit first ordering throughout the paper.

finite field (e.g. \mathbb{F}_p for some prime p). One notable feature of SPDZ is its 2-phase design: there is a *pre-processing phase* that produces correlated random values that are independent of the task to be securely computed, and the pre-computed random values will then be consumed in the *online phase* to enable very efficient computation. SPDZ aims to provide highly efficient online phase primitives such as secure addition and secure multiplication. Then high-level protocols can be implemented on top of SPDZ by calling the online phase primitives to compute a task expressed as an arithmetic circuit. In addition to efficiency, another benefit that SPDZ offers is strong security: it is UC secure against a static, active adversary corrupting up to $n - 1$ parties, and this strong security extends to high level protocols implemented on top of it.

On the technical side, SPDZ utilizes authenticated shares. In SPDZ, a value $x \in Z_p$ in the shared form is defined as:

$$[[x]] = (x_1, \dots, x_n, m_1^{(x)}, \dots, m_n^{(x)}, \Delta_1, \dots, \Delta_n),$$

and each party P_i holds a tuple $[[x]]_i = (x_i, m_i^{(x)}, \Delta_i)$ such that:

$$x = \sum_{i=1}^n x_i, \quad m^{(x)} = \sum_{i=1}^n m_i^{(x)}, \quad \Delta = \sum_{i=1}^n \Delta_i.$$

Each value is authenticated by a MAC. In the above, Δ is a global MAC key and the MAC is $m^{(x)} = x \cdot \Delta$. The authenticity of x can be verified by letting each P_i compute $\sigma_i = m_i^{(x)} - x \cdot \Delta_i$ and broadcast σ_i , then check if $\sum_{i=1}^n \sigma_i = 0$. The three parts in the tuple $[[x]]_i$ are additive shares of x , the MAC and the MAC key respectively.

In our protocols, we will explicitly use the following online phase primitives from SPDZ:

- $[[x + y]] \leftarrow [[x]] + [[y]]$: given shared values $[[x]]$ and $[[y]]$, compute the sum. This is done locally by each party P_i by computing $[[x + y]]_i = (x_i + y_i, m_i^{(x)} + m_i^{(y)}, \Delta_i)$.
- $[[a + x]] \leftarrow a + [[x]]$: add a shared value $[[x]]$ with a public value a . To do so, P_1 computes $[[a + x]]_1 = (x_1 + a, m_1 + a \cdot \Delta_1, \Delta_1)$, and each other party P_i computes $[[a + x]]_i = (x_i, m_i + a \cdot \Delta_i, \Delta_i)$.
- $[[a \cdot x]] \leftarrow a \cdot [[x]]$: multiply a shared value $[[x]]$ with a public value a . Each P_i computes locally $[[a \cdot x]]_i = (a \cdot x_i, a \cdot m_i, \Delta_i)$ from $[[x]]_i$.
- $\text{reveal}([[x]])$: reveal x in a shared value $[[x]]$, each P_i broadcasts x_i in $[[x]]_i$ and computes $x = \sum_{i=1}^n x_i$.
- $[[x \cdot y]] \leftarrow [[x]] \cdot [[y]]$: multiply two shared values. It is done by using Beaver's triple [10], i.e. a triple $(([a], [b], [c]))$ where a, b are random numbers in \mathbb{F}_p and $c = a \cdot b$. The triples are generated in the pre-processing phase. In the online phase when computing multiplication, a fresh random triple is used. It works by revealing (which requires broadcast) $[[\epsilon]]$ and $[[\rho]]$ where $[[\epsilon]] \leftarrow [[x]] - [[a]]$ and $[[\rho]] \leftarrow [[y]] - [[b]]$. Then

the product can be obtained as $[[x \cdot y]] \leftarrow [[c]] + \epsilon[[b]] + \rho[[a]] + \epsilon\rho$.

- $\text{Output}([[x]])$: this is used at the end of a protocol to output the final result x . It first checks the MACs of all values previously revealed in the protocol. If it fails, then aborts. Otherwise, it reveals x in $[[x]]$ to all parties, and checks the MAC of x . It aborts if it fails, and it outputs x otherwise.

Our protocols will use the pre-processing protocols in SPDZ for generating Beaver's triples. Since pre-processing is necessary for our protocols, we will treat the pre-processing phase as in place implicitly and not explicitly mention calling it, in the description of the protocols.

In SPDZ (and in many other secret-sharing MPC schemes), since computation over shares is simple modular addition and multiplication in a small finite field, the performance bottleneck of online protocols is often network communication [40, 41]. Therefore, reducing the number of rounds and number of interactions is crucial to the efficiency of the online protocols.

3.3 Differential Privacy

Differential privacy [24] is a well-established principle that quantifies the privacy impact on individuals, when their private information is included in a dataset and some statistics obtained from the dataset are released. The first definition of differential privacy is the following:

Definition 1 (ϵ -differential privacy [24]). *A randomized mechanism $f : \mathcal{D} \rightarrow \mathcal{R}$ gives ϵ -differential privacy, where ϵ is a positive real number, if for all data sets D_1 and D_2 differing in at most one element, and all $R \subseteq \mathcal{R}$,*

$$e^{-\epsilon} \cdot \Pr[f(D_2) \in R] \leq \Pr[f(D_1) \in R] \leq e^{\epsilon} \cdot \Pr[f(D_2) \in R].$$

Definition 1 is very strong but also often renders the output unusable, since it incurs substantial distortion to be enforced. Therefore, (ϵ, δ) -differential privacy is often used:

Definition 2 ((ϵ, δ) -differential privacy [25]). *A randomized mechanism $f : \mathcal{D} \rightarrow \mathcal{R}$ gives (ϵ, δ) -differential privacy, where (ϵ, δ) are positive real numbers, if for all data sets D_1 and D_2 differing in at most one element, and all $R \subseteq \mathcal{R}$,*

$$e^{-\epsilon} \cdot \Pr[f(D_2) \in R] - \frac{\delta}{e^{\epsilon}} \leq \Pr[f(D_1) \in R] \leq e^{\epsilon} \cdot \Pr[f(D_2) \in R] + \delta.$$

Intuitively, (ϵ, δ) -differential privacy ensures that for all adjacent D_1, D_2 , the absolute value of the privacy loss will be bounded by ϵ with a probability at least $1 - \delta$.

3.4 Statistical Security

We briefly review the notion of statistical security [32] that we use in our ZeroTest sub-protocol (see Section 4.5). This

notion requires that the views of protocol execution can be simulated such that the distributions of real and simulated views are statistically indistinguishable. Formally, let X and Y be distributions with finite sample spaces V and W and $\Delta(X, Y) = \frac{1}{2} \sum_{v \in V \cup W} |Pr(X = v) - Pr(Y = v)|$ the statistical distance between them. We say that the distributions are statistically indistinguishable if $\Delta(X, Y) \leq \text{negl}(\lambda)$ where negl is a negligible function and λ is some statistical security parameter. As usual, a function is negligible if for every positive polynomial p there is an N such that for all integers $n > N$ it holds that $\text{negl}(n) < \frac{1}{p(n)}$. Statistical security is information theoretic, i.e. it holds even if the adversary has unbounded computational power. The statistical security parameter usually can be smaller than the computational security parameter (e.g. 40 is often used in the literature [36, 50]).

3.5 Universal Composability (UC)

We briefly review the UC framework [12] that we use to prove the security of our protocol. Being UC secure means that our protocol can be freely composed with other protocols and still be secure. The UC framework is defined in terms of comparing a real world execution and the execution in an ideal world, in the presence of an adversary (environment). Security in UC is defined in terms of the adversary's inability to distinguish whether it is interacting with the real protocol Π , or with a simulator in the ideal world which has access to an ideal functionality \mathcal{F} . If so, then we say that the protocol Π securely realizes the functionality \mathcal{F} . Intuitively, the ideal world is secure by definition, and a successful simulation means that the adversary running the protocol in real world cannot do more damage than what is allowed in the ideal world, hence the protocol is secure.

Let the adversary be \mathcal{Z} . In the beginning of an execution, \mathcal{Z} chooses inputs for all parties and gets their outputs when the execution finishes. It also controls some corrupted parties, which means \mathcal{Z} will instruct what they should do during the execution and see the communication and internal states of them. When \mathcal{Z} stops, it outputs a bit. Security is established by showing the existence of a simulator \mathcal{S} that interacts with both \mathcal{F} and \mathcal{Z} . The simulator should be able to simulate the view of the protocol that looks like what \mathcal{Z} would see in a real attack by playing the honest parties' role when interacting with \mathcal{Z} , but without access to the input and state of the honest parties. One significant difference in the simulation in UC and in stand-alone environment is that \mathcal{Z} can query the corrupted parties during the execution (rather than just collect the views after the execution). This means some techniques such as rewinding cannot be used in UC proofs. For a more formal and complete account of the UC framework, please refer to [12].

4 The PDCE Protocol

4.1 Overview

In the PDCE protocol, we have a set of n honest Data Parties (DPs) and a set of d untrusted (up to $d - 1$ can be malicious) Computation Parties (CPs). The DPs are responsible for data collection. They observe the events of interest, e.g. IP addresses of the visitors, and record them locally as a set of FM sketches. After the data has been collected, the DPs secret-share the sketches among the CPs, who will securely combine them, and compute the estimator Z_N of the count of distinct values. The protocol has four phases: initialization phase, offline phase, data collection phase, and data aggregation phase. Each phase involves certain sub-protocols.

4.2 Initialization Phase

In this phase, the parties negotiate parameters to be used in the protocol. This phase only needs to run once when setting up the system. Firstly, all parties need to agree on a finite field \mathbb{F}_p . This field will be used as the basis of data representation, secret sharing and all computation. The modulus p is decided by three parameters: (1) λ , which is a statistical security parameter (e.g. 40); (2) τ , which determines the size of the plaintext domain (integers between $[0, 2^\tau - 1]$); (3) M , e.g. 32768, which comes from the BGV somewhat homomorphic encryption [11] used by SPDZ. Specifically, the parties choose p that is a $(\lambda + \tau)$ -bit prime number and M divides $p - 1$. Next, the parties agree on the parameters for FM sketches. Given the accuracy and privacy requirements, they decide m (the number of sketches to be used). Based on the pre-knowledge of the maximum number of items that can be observed collectively, the parties decide w (the size of each sketch). Finally, the CPs run the setup protocol of SPDZ to obtain the parameters and keys for SPDZ.

4.3 Offline Phase

In the offline phase, the CPs run the pre-processing protocol of SPDZ. In addition, they also run a few other offline protocols to generate various random values that will be used later in the data collection and aggregation phases. The offline protocols we use already exist in the literature, therefore we only give a high level description of them here. The protocol details and references can be found in the full version.

- $\text{Rand}()$: generates $\llbracket r \rrbracket$, the shares of a random value $r \in_R \mathbb{F}_p$.
- $\text{Rand2}()$: generates $\llbracket b \rrbracket$, the shares of a random bit $b \in_R \{0, 1\}$.
- $\text{RandExp}(l)$: generates $(\llbracket R^{-1} \rrbracket, \llbracket R \rrbracket, \llbracket R^2 \rrbracket, \dots, \llbracket R^l \rrbracket)$, the shares of a random number $R \in_R \mathbb{Z}_p^*$, as well as the shares of its i th powers (for $i = -1$ and $2 \leq i \leq l$).

4.4 Data Collection Phase

At the beginning of this phase, the DPs choose a keyed hash function H , a pseudorandom function PRF , and establish a secret key sk for PRF among them. The secret key sk can be established using an authenticated group key exchange protocol (e.g. [39]). The PRF and the key sk will be used for deriving hash keys, so that m independent FM sketches can be constructed using H and different hash keys. For $1 \leq j \leq m$, the j th hash key is $k_j = PRF(sk, j)$. Then each DP maintains m FM sketches, observes items and adds them into its FM sketches. At the end of this phase, each DP splits its FM sketches into secret shared form, and sends the shares to the CPs. The protocol for data collection is shown in Protocol 1, and the sub-protocol $Share(x)$ is shown in Protocol 2.

Protocol 1: Data Collection

Input: Each DP's input is sk , the shared key for the PRF
Result: The CPs obtain the shares of the FM sketches

```
// Initialize FM sketches
1 Each  $DP_i$  initialize  $m$  FM sketches, each is  $w$ -bit
// Collect data
2 Whenever  $DP_i$  observes an item  $o$ , it does the following:
// add  $o$  to sketch (see Sec. 3.1)
3   for  $j = 1; j \leq m; j++$  do
4     Compute  $l = \rho(H(k_j || o))$ ;
5     Set  $FS_i^j[l] = 1$ ;
6   end
// Finish data collection
7 After data has been collected, each  $DP_i$  does the following:
8   for  $j = 1; j \leq m; j++$  do
9     for  $l = 0; l \leq w - 1; l++$  do
10      Run  $Share(FS_i^j[l])$  with the CPs;
11    end
12  end
```

Protocol 2: $Share(x)$

Offline: CPs run $\llbracket a \rrbracket \leftarrow Rand()$, where $a \in_R \mathbb{F}_p$.
Input: The DP's input is x , the value to be shared.
Result: The CPs obtain $\llbracket x \rrbracket$

```
1 CPs reveal  $a$  to DP;
2 DP computes  $x - a$  and broadcasts it to all CPs;
3 CPs obtain  $\llbracket x \rrbracket = \llbracket a \rrbracket + (x - a)$ ;
```

4.5 Data Aggregation Phase

This phase involves only the CPs. The CPs first merge the shares from the DPs into m shared FM sketches such that each slot in the sketches holds either a zero or a positive integer. Then, they convert the integer FM sketches into binary FM sketches. After that, they extract the estimator Z_N from the sketches, and compute the count from the estimator locally.

Merge Shares At the start of the data aggregation phase, each CP holds the shares of all the FM sketches from all DPs. The first step for each CP is to merge the shares of

the sketches to get the shares of a set of m (integer) FM sketches that record the union of observations from all DPs. As mentioned in Section 3.1, merging FM sketches can be done by bit-wisely ORing the sketches. However, the Boolean OR operation corresponds to multiplication of shared values. A naive implementation of this step thus would require $(n - 1) \cdot m \cdot w$ multiplication operations and thus $(n - 1) \cdot m \cdot w$ rounds of communication, where n is the number of DPs, m is the number of FM sketches generated by each DP, and w is the bit-size of the FM sketches. To reduce the cost, in our protocol, we merge the shares by addition. The protocol is shown in Protocol 3. For the l -th bit in the j -th FM sketches, the CPs locally sum up the n shares for that bit from all DPs. At the end, the CPs obtain the shares of m integer FM sketches such that 0 in the integer FM sketches corresponds to 0 in binary FM sketches, and non-zero corresponds to 1. The integer FM sketches will be converted to binary sketches in the next step. The only operation needed in this step is addition. Thus, no interaction is required. Looking ahead, the next step requires in total $2 \cdot m \cdot w$ rounds of interaction, thus the total cost is much less than the naive implementation in real applications where the number of DP is often large.

Protocol 3: MergeShares

Input: Each CP_k holds $\llbracket FS_i^j[l] \rrbracket_k$
 $(1 \leq i \leq n, 1 \leq j \leq m, 0 \leq l \leq w - 1)$
Result: $\llbracket FS_{\cup}^j[l] \rrbracket_k$ ($1 \leq j \leq m, 0 \leq l \leq w - 1$)

```
1 for  $j = 1; j \leq m; j++$  do
2   for  $l = 0; l \leq w - 1; l++$  do
3      $\llbracket FS_{\cup}^j[l] \rrbracket_k = \sum_{i=1}^n (\llbracket FS_i^j[l] \rrbracket_k)$ ;
4   end
5 end
```

Protocol 4: $ToBinary(\llbracket FS_{\cup}^1[0] \rrbracket, \dots, \llbracket FS_{\cup}^1[w - 1] \rrbracket, \dots, \llbracket FS_{\cup}^m[0] \rrbracket, \dots, \llbracket FS_{\cup}^m[w - 1] \rrbracket)$

Input: $\llbracket FS_{\cup}^j[l] \rrbracket$ ($1 \leq j \leq m, 0 \leq l \leq w - 1$), shares of the m integer FM sketches.
Result: $\llbracket BFS_{\cup}^j[l] \rrbracket$ ($1 \leq j \leq m, 0 \leq l \leq w - 1$), shares of the m converted binary FM sketches.

```
1 for  $j = 1; j \leq m; j++$  do
2   for  $l = 0; l \leq w - 1; l++$  do
3      $\llbracket BFS_{\cup}^j[l] \rrbracket = ZeroTest(\llbracket FS_{\cup}^j[l] \rrbracket)$ ;
4   end
5 end
```

Convert to Binary Sketches As shown in Protocol 4, the second step is to covert each FS_{\cup}^j back to the normal binary FM sketches², so that we can later extract the estimator from them. This is done by running a zero test protocol among the CPs on each slot that sets the slot to 0 if the value stored in it is 0, or to 1 otherwise.

Here we use the protocol from [42]. The protocol is based on the following idea: to test whether a is 0 or not, we first

²To clarify, here binary means $\{0, 1\}$ in \mathbb{F}_p , not $\{0, 1\}$ in \mathbb{F}_2

Protocol 5: ZeroTest($\llbracket a \rrbracket$)

Offline:**for** $i = 0, \dots, l-2$, where l is the bit length of p **do**
 $\llbracket r_i \rrbracket \leftarrow \text{Rand2}()$;**end** $\llbracket r \rrbracket \leftarrow \sum_{i=0}^{l-2} 2^{l-2-i} \llbracket r_i \rrbracket$;

// interpolate the lookup polynomial

 $(\tau, \beta_0, \dots, \beta_\tau) \leftarrow \text{interpolate}()$;**Input:** $\llbracket a \rrbracket$, where a is a τ -bit integer.**Result:** $\llbracket b \rrbracket$, where $b = 0$ if $a = 0$, $b = 1$ otherwise

- 1 $\llbracket m \rrbracket = \llbracket r \rrbracket + \llbracket a \rrbracket$;
 - 2 Reveal $\llbracket m \rrbracket$;
 - 3 $\llbracket 1+h \rrbracket = 1 + \sum_{i=l-1}^{\tau} (\llbracket r_i \rrbracket + m_i - 2\llbracket r_i \rrbracket \cdot m_i)$;
 - 4 $\llbracket b \rrbracket = \text{Lookup}(\llbracket 1+h \rrbracket, \tau, \beta_0, \dots, \beta_\tau)$
-

Protocol 6: Lookup($\llbracket x \rrbracket, \ell, \beta_0, \dots, \beta_\ell$)

Offline: $(\llbracket R^{-1} \rrbracket, \llbracket R \rrbracket, \llbracket R^2 \rrbracket, \dots, \llbracket R^\ell \rrbracket) \leftarrow \text{RandExp}(\ell)$;**Input:** $\llbracket x \rrbracket$, where x is an integer; ℓ is the degree of the lookup polynomial $f(\cdot)$; $\beta_0, \dots, \beta_\ell$ are the coefficient of $f(\cdot)$.**Result:** $\llbracket y \rrbracket$, where $y = f(x)$.

- 1 $\llbracket a \rrbracket = \llbracket R^{-1} \rrbracket \cdot \llbracket x \rrbracket$;
 - 2 Reveal $\llbracket a \rrbracket$;
 - 3 **for** $i = 2, \dots, \ell$ **do**
 - 4 $\llbracket x^i \rrbracket = a^i \cdot \llbracket R^i \rrbracket$
 - 5 **end**
 - 6 $\llbracket b \rrbracket = \sum_{i=0}^{\ell} \beta_i \cdot \llbracket x^i \rrbracket$
-

compute $r + a$ where r is a random integer, and then compute the Hamming distance h between $r + a$ and r . Obviously, if $a = 0$, then $h = 0$; otherwise h is a small integer in $[1, \tau]$, where τ is the bit length of the plaintext. As h is small, it is feasible to use a lookup function that is a polynomial $f(\cdot)$ such that $f(0) = 0$ and $f(x) = 1$ for all other $x \in [1, \tau]$. There is a small technicality that $f(0)$ cannot be evaluated without leaking information. To see that, note that in the first line of Protocol 6, if $x = 0$ then $a = 0$, and revealing a will reveal whether x is 0. Thus in line 3 of Protocol 5, 1 is added to h so that the input to the polynomial will never be 0. The lookup polynomial will be interpolated accordingly (e.g. using Lagrange Interpolation), and evaluating f at $h + 1$ will output 0 if h is 0 or 1 otherwise. The *ZeroTest* protocol is shown in Protocol 5, and the sub-protocol *Lookup* for evaluating the lookup function is shown in Protocol 6 (both are from [42]).

Extract Estimator Recall that given an FM sketch, one can extract z_N , i.e. the index of the first 0 bit in the sketch. When using m FM sketches, the sum $Z_N = \sum_{i=1}^m z_{N,i}$ will be used to estimate the number of distinct observed items as $\tilde{N} = \frac{2^{\frac{Z_N}{m}} - 2^{-k \cdot \frac{Z_N}{m}}}{\phi}$ (see Section 3.1). The formula is deterministic and invertible, therefore revealing \tilde{N} and revealing Z_N are essentially equivalent. Because of this, we can let the protocol output Z_N rather than \tilde{N} without compromising correctness or security. With Z_N , each CP can locally compute \tilde{N} .

Z_N can be extracted using the following simple idea: firstly,

Protocol 7: ExtractZ($\llbracket \text{BFS}_{\cup}^1[0] \rrbracket, \dots, \llbracket \text{BFS}_{\cup}^1[w-1] \rrbracket, \dots, \llbracket \text{BFS}_{\cup}^m[0] \rrbracket, \dots, \llbracket \text{BFS}_{\cup}^m[w-1] \rrbracket$)

Input: $\llbracket \text{BFS}_{\cup}^1[0] \rrbracket, \dots, \llbracket \text{BFS}_{\cup}^1[w-1] \rrbracket, \dots, \llbracket \text{BFS}_{\cup}^m[0] \rrbracket, \dots, \llbracket \text{BFS}_{\cup}^m[w-1] \rrbracket$, the shares of the m binary FM sketches.**Result:** Z_N , the estimator extracted from the sketches

- 1 $\llbracket Z_N \rrbracket = 0$;
 - 2 **for** $i = 1; i \leq m; i++$ **do**
 - 3 $\llbracket Z_N \rrbracket = \llbracket Z_N \rrbracket + \llbracket \text{BFS}_{\cup}^i[0] \rrbracket$;
 - 4 **end**
 - 5 **for** $l = 1; l \leq w-1; l++$ **do**
 - 6 **for** $i = 1; i \leq m; i++$ **do**
 - 7 $\llbracket \text{BFS}_{\cup}^i[l] \rrbracket = \llbracket \text{BFS}_{\cup}^i[l-1] \rrbracket \cdot \llbracket \text{BFS}_{\cup}^i[l] \rrbracket$;
 - 8 $\llbracket Z_N \rrbracket = \llbracket Z_N \rrbracket + \llbracket \text{BFS}_{\cup}^i[l] \rrbracket$;
 - 9 **end**
 - 10 **end**
 - 11 return $Z_N \leftarrow \text{Output}(\llbracket Z_N \rrbracket)$;
-

for each sketch, set all bits after the first 0 bit to 0, then sum up all bits in the sketch, and the result is the estimator $z_{N,i}$; then Z_N can be obtained by summing up all $z_{N,i}$'s. This is essentially what we do in Protocol 7. To set bits to 0 after the first 0 bit, the protocol does the following: for each sketch FS_{\cup}^i , it sets $\text{FS}_{\cup}^i[l] = \text{FS}_{\cup}^i[l-1] \cdot \text{FS}_{\cup}^i[l]$ sequentially for $1 \leq l \leq w-1$. By doing so, all bits before and including the first 0 bit remain unchanged, and all bits after will be set to 0 due to the chained multiplication.

Although Protocol 7 is simple, it requires $w-1$ rounds because the multiplication in each round is dependent on the output of the previous round. To improve the round efficiency, we designed another protocol (Protocol 10) that only requires $2 \log(w)$ rounds. The protocol can be found in Appendix A. As we will see later in Section 6, the performance of Protocol 10 is better than Protocol 7 when the network bandwidth is limited.

Estimate the Cardinality After extracting Z_N , each CP computes the estimated cardinality locally with the formula $\tilde{N} = \frac{2^{\frac{Z_N}{m}} - 2^{-k \cdot \frac{Z_N}{m}}}{\phi}$, as explained in Section 3.1.

5 Security Analysis

5.1 Protocol Security

We prove the security of the protocol in Section 4 in the Universally Composable framework [13]. This provides a strong notion of security and allows our protocol to serve as a component of a larger system without losing its security properties. We adopt a very strong adversary model in which the untrusted CPs are modelled as corrupted by a single adversary that is malicious, i.e. can behave arbitrarily. The adversary can corrupt all but one CPs statically. Informally, with only one honest CP, the security properties of the protocol are: (1) the adversary learns nothing from executing the protocol except the differentially private output of the protocol; (2) the adver-

sary cannot affect the correctness of the computation without being detected. The security properties we considered in this paper are confidentiality and correctness. We leave out other properties such as robustness. In essence, the protocol will terminate if any party aborts and no result will be computed. This limitation is inherent in the underlying MPC framework we use, namely SPDZ. That said, robust MPC is an active research topic and our protocol can be migrated to a robust MPC framework when it is available.

Our adversary model is quite similar to that in [30], except that (1) in our model, the DPs are honest but in [30] they allow DPs to be corrupted adaptively; (2) In our model, malicious CPs cannot tamper with the data and the result, while in [30] a malicious CP can insert elements into the hashtable and change the result (and this cannot be prevented unless their protocol is significantly changed). Regarding whether the DPs should be assumed honest or not, we have the following remarks: (1) We model the DPs as honest mainly because, like many other differential privacy mechanisms, we need to keep the randomness, namely the PRF key, private from the adversary. Compromising this key will break the differential privacy guarantee. On the other hand, although [30] allows DPs to be corrupted, once a DP is corrupted, differential privacy guarantee is broken as well. This is because the adversary can now see the raw data collected by the corrupted DP. If the element x that differentiates D_1 and D_2 happens to be observed by the adversary, differential privacy is broken. (2) After corrupting a DP, [30] can prevent the adversary from seeing the corrupted DPs' data before corruption. This is something our protocol cannot achieve now. However, firstly [30] is used for Tor, and they consider law enforcement forcing DPs to reveal data collected a threat, but this is not common in other applications; secondly, we can easily achieve it, by secret-sharing the FM sketches when initializing them, and update them obliviously. This only adds one round of communication between DPs and CPs, and negligible computation. (3) Our DP side computation is cheap (hashing) and requires only small storage (a few MB for thousands of FM sketches and one secret key). Thus, it is relatively easy to secure DPs, e.g. using trusted hardware like Intel SGX. Spending reasonable efforts on securing DPs in exchange for much less computation on CPs seems to be a worthy trade-off.

Note, as in many proofs, we prove the security modularly in the so called \mathcal{F} -hybrid model. That is, we can replace an already proven secure sub-protocol with an ideal functionality. Theorem 5 states two ideal functionalities; $\mathcal{F}_{\text{SPDZ}}$ and $\mathcal{F}_{\text{offline}}$. The first is the ideal functionality for the SPDZ protocol, whose security has been proven in [19, 41]. The second is the ideal functionality for our offline protocols. The offline protocols are from the literature, therefore we also separate them as an ideal functionality. The details of $\mathcal{F}_{\text{offline}}$ as well as the full security proof (under the SPDZ framework) can be found in the full version. Then, the security properties of the online protocol that does the cardinality estimation are

captured by an ideal functionality in Figure 1. We have the following theorem:

Theorem 1. *In the $\mathcal{F}_{\text{SPDZ}}, \mathcal{F}_{\text{offline}}$ -hybrid model, the protocol in Section 4 realizes $\mathcal{F}_{\text{PDCE}}$ with statistical security against any malicious adversary who statically corrupts up to $d - 1$ CPs.*

The proof of Theorem 1 can be found in the full version.

Functionality $\mathcal{F}_{\text{PDCE}}$
The functionality maintains a dictionary, Val , to keep track of the authenticated values. Entries of Val lie in the (fixed) finite field \mathbb{F}_p and cannot be changed, for simplicity.
Abort: On receiving Abort from the adversary, send \perp to all parties and terminate.
Share: On receiving (share, x, id) from DP , and (share, id) from all CPs, set $\text{Val}[id] \leftarrow x$.
Go: After receiving (go) from all parties, ignore messages from DP and the following methods can be called from now on.
MergeShare: On receiving $(\text{mergeshare}, id^{\text{FS}}, id^{\text{FS}_U})$ from all CPs, where id^{FS} is a $(mw \times n)$ matrix and id^{FS_U} is an mw vector, all contain some ids, set for $1 \leq i \leq mw$, $\text{Val}[id_i^{\text{FS}_U}] \leftarrow \sum_{j=1}^n \text{Val}[id_{i,j}^{\text{FS}}]$.
Lookup: On receiving $(\text{lookup}, id_x, id_y, \ell, \beta_0, \dots, \beta_\ell)$ from all CPs, check that $\ell, \beta_0, \dots, \beta_\ell$ defines a lookup polynomial as expected, then set $\text{Val}[id_y] \leftarrow \sum_{i=0}^{\ell} \beta_i \cdot (\text{Val}[id_x])^i$.
ZeroTest: On receiving $(\text{zerotest}, id_a, id_b)$ from all CPs, if $\text{Val}[id_a] = 0$, set $\text{Val}[id_b] \leftarrow 0$, otherwise set $\text{Val}[id_b] \leftarrow 1$.
ExtractZ: On receiving $(\text{extractZ}, id_0^1, \dots, id_{w-1}^1, id_0^2, \dots, id_{w-1}^2, \dots, id_0^m, \dots, id_{w-1}^m, id_{Z_N})$ from all CPs, count from the beginning the number of continuous 1 in $(\text{Val}[id_0^1], \dots, \text{Val}[id_{w-1}^1])$ to get $z_{N,i}$, then compute $Z_N = \sum_{i=1}^m z_{N,i}$, set $\text{Val}[id_{Z_N}] \leftarrow Z_N$.

Figure 1: Ideal Functionality for the PDCE Protocol

5.2 Differential Privacy

In this section, we will show that if the cardinality to be estimated by the FM sketches is large enough (larger than a threshold N_0), then our protocol in Section 4 satisfies (ϵ, δ) -differential privacy automatically, without requiring any further manipulation of the output. We noticed that in [21], the authors conclude that cardinality estimation by sketches does not preserve privacy. However, our positive result does not contradict their negative result. The reason is that in their model, the adversary can access the sketches and the final estimation result; while in our model, since MPC is used, the adversary can only access the final estimation result. The sketches are secret-shared in the protocol and are never revealed (if at least one CPs is honest). In fact, the mitigation strategies proposed in [21] are about restricting the adversary's access to the sketches, which is in line with what we do.

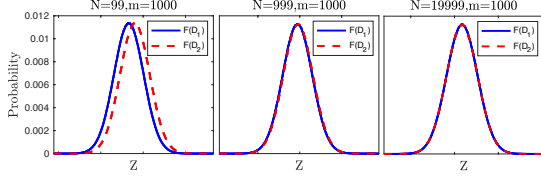


Figure 2: Results from the Monte Carlo Simulations

5.2.1 Intuition

The notion of (ϵ, δ) -differential privacy requires that the outputs from a randomized mechanism on two neighboring datasets should be close enough with high probability.

To start with, our protocol can be viewed as a randomized mechanism $F : 2^U \rightarrow \mathbb{Z}$. Let U be the universe of elements and $D \subseteq U$ a set comprised of the union of the observations of all n DPs. F takes D as input, internally builds m FM sketches of D , and outputs the random variable $Z_N = z_{N,1} + \dots + z_{N,m}$. Note that F is a randomized mechanism because a randomly chosen key is used in our protocol and is renegotiated in each run. The random key is used to derive the hash keys for each FM sketch. The use of hash keys also ensures the independence of $z_{N,i}$'s, albeit they might be generated on correlated data. To see that, for H that is modeled as a perfect random function, $\Pr[H(k_j|x) = y]$ is uniform and independent of $\Pr[H(k_j|x) = y']$. The independence of the hash output then implies $z_{N,i}$'s are independent. Also, as we mentioned in Section 3.1, FM sketches can eliminate duplicates in data because the same element will end up with the same hash value when hashed under the same hash key, thus multiple copies of the same element will be counted as one. Although data is collected by individual DPs, we can think the final result is about the union set of the elements from all DPs. We can model F just with one input D that is the union of the n sets from the DPs.

The output Z_N from F is a random variable which can take larger values as the cardinality of D increases. Intuitively, as D becomes larger, each element in D has a smaller contribution to Z_N . Eventually, the contribution becomes so insignificant and each element's presence will have almost no effect on the distribution of Z_N . In other words, when D is large enough, the addition or removal of an element from D will cause almost no change to the distribution of Z_N , i.e., differential privacy can be achieved. To illustrate the intuition, we conducted three Monte Carlo simulations. The results are shown in Fig. 2. In each simulation, two sets D_1 and D_2 were used, such that D_1 had N elements and D_2 was obtained from D_1 by adding one extra element. We set $N = 99$, $N = 999$ and $N = 19999$ in the three simulations. Each simulation had 10 million rounds. In each round, we generated a random set of m hash keys, built m sketches for D_1 and m sketches for D_2 , and then computed $F(D_i)$ from the sketches. Fig. 2 shows the distributions of $F(D_1)$ and $F(D_2)$, each obtained from 10 million samples. As can be seen, when N becomes larger, the two curves become closer.

Note that our protocol is designed for applications that

require one-off or periodical release of statistics (e.g. the number of distinct IP addresses per hour). In each run of our protocol, fresh randomness is introduced by renegotiating the PRF key, so that the sketches are independent of those in the previous run. The protocol does not use sketches from previous runs, and only one query is answered in each run (i.e. the output of each run is Z_N). The protocol does not support correlated queries, e.g. how many new elements have been added since the last estimation. If our protocol is used for answering correlated queries, differential privacy may no longer hold because correlated queries leak more information.

In the following, we will start by showing that when using a single FM sketch, we can find an N_0 such that the protocol satisfies (ϵ, δ) -differential privacy whenever the input set to the protocol has cardinality at least N_0 . Then the bound N_0 for (ϵ, δ) -differential privacy to hold in the m FM sketches case can be obtained by using the composition theorems of differential privacy [26]. The bound obtained from the composition theorems can be refined, to get a much smaller (better) N_0 .

5.2.2 Finding N_0 : Single FM Sketch³

Let z_N denote the discrete random variable extracted from an FM sketch when the input cardinality is N . We first work out the probability mass function (PMF) of z_N . In [31], the complementary cumulative distribution function of z_N was given as:

$$q_{N,k} = \Pr(z_N \geq k) = \sum_{j=0}^{2^k} (-1)^{v(j)} e^{-\frac{jN}{2^k}}$$

where $0 \leq k \leq w-1$ and $v(j)$ denotes the number of ones in the binary representation of j . Then, we can derive the PMF of z_N as:

$$p_{N,k} = \Pr(z_N = k) = q_{N,k} - q_{N,k+1}$$

The above, after some derivation, gives us:

$$p_{N,k} = \begin{cases} e^{-\frac{N}{2}} & \text{if } k = 0 \\ e^{-\frac{N}{2^{k+1}}} \prod_{j=0}^{k-1} (1 - e^{-\frac{N}{2^{j+1}}}) & \text{if } k > 0 \end{cases} \quad (2)$$

We want (ϵ, δ) -differential privacy to hold for any sufficiently large datasets D_1 and D_2 differing in at most one element. When using a single FM sketch in our protocol, it is equivalent to say that we want to find an N_0 such that for all $N \geq N_0$ and for all k , the following holds:

$$\begin{cases} \Pr[z_N = k] \leq e^\epsilon \cdot \Pr[z_{N+1} = k] + \delta \\ \Pr[z_{N+1} = k] \leq e^\epsilon \cdot \Pr[z_N = k] + \delta \end{cases}$$

which is equivalent to:

$$e^{-\epsilon} \cdot p_{N+1,k} - \frac{\delta}{e^\epsilon} \leq p_{N,k} \leq e^\epsilon \cdot p_{N+1,k} + \delta.$$

It is easy to see that the above holds, if at each k either of the following two conditions is true: (1) $e^{-\epsilon} \leq \frac{p_{N,k}}{p_{N+1,k}} \leq e^\epsilon$

³More details and proofs can be found in the full version.

(ϵ -differential privacy holds at those k), or (2) $p_{N,k} \leq \delta$ and $p_{N+1,k} \leq \delta$ (the probability of getting to this k is sufficiently small). When condition (1) is true, (ϵ, δ)-differential privacy holds because

$$e^{-\epsilon} \cdot p_{N+1,k} - \frac{\delta}{e^\epsilon} < e^{-\epsilon} \cdot p_{N+1,k} \leq p_{N,k} \leq e^\epsilon \cdot p_{N+1,k} < e^\epsilon \cdot p_{N+1,k} + \delta$$

When condition (2) is true, (ϵ, δ)-differential privacy holds because

$$e^{-\epsilon} \cdot p_{N+1,k} - \frac{\delta}{e^\epsilon} < 0 \leq p_{N,k} \leq \delta \leq e^\epsilon \cdot p_{N+1,k} + \delta$$

To start with, we prove the following lemma:

Lemma 1. $\frac{p_{N,k}}{p_{N+1,k}}$ decreases monotonically in k .

Looking ahead, based on Lemma 1, our strategy for finding N_0 consists of two steps:

1. Find N_1 such that for all $N \geq N_1$, there exists k_{min} , and (i) for all $k \geq k_{min}$, $\frac{p_{N,k}}{p_{N+1,k}} \leq e^\epsilon$ and (ii) for all $k < k_{min}$, $p_{N,k} \leq \delta$ and $p_{N+1,k} \leq \delta$.
2. Find N_2 such that for all $N \geq N_2$, there exists k_{max} , and (i) for all $k \leq k_{max}$, $\frac{p_{N,k}}{p_{N+1,k}} \geq e^{-\epsilon}$ and (ii) for all $k > k_{max}$, $p_{N,k} \leq \delta$ and $p_{N+1,k} \leq \delta$.

Then, we take $N_0 = \max(N_1, N_2)$. Clearly, for all $N \geq N_0$, we have: (i) $e^{-\epsilon} \leq \frac{p_{N,k}}{p_{N+1,k}} \leq e^\epsilon$ for $k_{min} \leq k \leq k_{max}$, and (ii) $p_{N,k} \leq \delta$ and $p_{N+1,k} \leq \delta$ for $k < k_{min}$ or $k > k_{max}$. Thus, (ϵ, δ)-differential privacy holds for all $N \geq N_0$ and all k (see Figure 3).

Probability

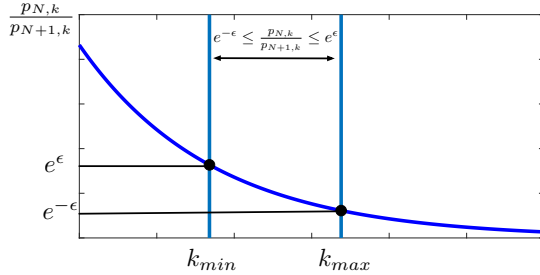
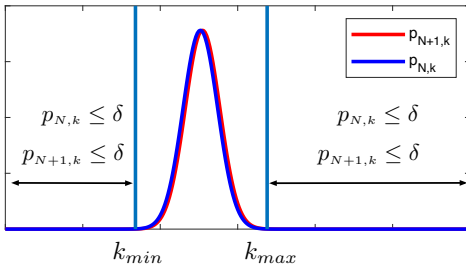


Figure 3: k_{min} and k_{max}

Finding k_{min} and N_1 We first show the existence of k_{min} :

Lemma 2. Let $k_{min} = \max(\lceil \log_2 \frac{1}{\epsilon} \rceil - 1, 0)$. For any $\epsilon > 0$ and any $k \geq k_{min}$, it holds that $\frac{p_{N,k}}{p_{N+1,k}} \leq e^\epsilon$.

Lemma 2 tells us that k_{min} always exists for any $\epsilon > 0$, and that it is independent of N . Next we will show that the increase of N can eventually make $p_{N,k} \leq \delta$ and $p_{N+1,k} \leq \delta$ for all $k < k_{min}$. First, we prove the following lemma:

Lemma 3. For all $k_{min} > 0$, $\Pr[z_N < k_{min}]$ decreases monotonically in N , and $\lim_{N \rightarrow \infty} \Pr[z_N < k_{min}] = 0$.

Now we are ready to state the following theorem:

Theorem 2. Let $k_{min} = \max(\lceil \log_2 \frac{1}{\epsilon} \rceil - 1, 0)$ and N_1 be the smallest positive integer such that $1 - q_{N_1, k_{min}} \leq \delta$. Then, for all $N \geq N_1$, it holds that $p_{N,k} \leq \delta$ and $p_{N+1,k} \leq \delta$ when $k < k_{min}$, and $\frac{p_{N,k}}{p_{N+1,k}} \leq e^\epsilon$ when $k \geq k_{min}$.

Finding k_{max} and N_2 We now show the existence of k_{max} . Note that unlike k_{min} , k_{max} is a value that is dependant on N .

Lemma 4. Let $k_{max} = \lceil \log_2 N \rceil + c$ and $c = \lceil \frac{-1 + \sqrt{1 + 8 \log_2 \frac{1}{\delta}}}{2} \rceil$. For all $0 < \delta < 1$ and $N \in \mathbb{Z}^+$, $p_{N,k} \leq \delta$ and $p_{N+1,k} \leq \delta$ for all $k > k_{max}$.

Next we want to find N_2 such that for all $N \geq N_2$, $\frac{p_{N, k_{max}}}{p_{N+1, k_{max}}} \geq e^{-\epsilon}$. If this holds, then by Lemma 1, $\frac{p_{N,k}}{p_{N+1,k}} \geq e^{-\epsilon}$ for all $k \leq k_{max}$. Recall that for $k = 0$, $\frac{p_{N,k}}{p_{N+1,k}} = e^{\frac{1}{2}} > e^{-\epsilon}$ for all N trivially. Then we only need to consider the case $k > 0$.

In this case, $\frac{p_{N,k}}{p_{N+1,k}} = e^{\frac{1}{2^{k+1}}} \cdot \prod_{j=0}^{k-1} \frac{(1 - e^{-\frac{N}{2^{j+1}}})}{(1 - e^{-\frac{N+1}{2^{j+1}}})}$. Let us define

$$\Psi(N, k) = \prod_{j=0}^{k-1} \frac{(1 - e^{-\frac{N}{2^{j+1}}})}{(1 - e^{-\frac{N+1}{2^{j+1}}})} \quad (3)$$

We can see if $\Psi(N, k) \geq e^{-\epsilon}$ then $\frac{p_{N,k}}{p_{N+1,k}} \geq e^{-\epsilon}$, because $e^{\frac{1}{2^{k+1}}} \geq 1$.

It is actually not difficult to find some N_2 such that $\Psi(N_2, k_{max}) \geq e^{-\epsilon}$. The tricky part is whether for all $N \geq N_2$, $\Psi(N, k_{max}) \geq e^{-\epsilon}$ still holds. If $\Psi(N, k_{max})$ is monotonically increasing in N , then this can be proved. However, this is only partially true. Regarding this, we have the following:

Lemma 5. Let k_{max} as defined in Lemma 4, $\Psi(N, k_{max}) < \Psi(N+1, k_{max})$ if $\lceil \log_2 N \rceil = \lceil \log_2(N+1) \rceil$.

In the case that $\lceil \log_2 N \rceil \neq \lceil \log_2(N+1) \rceil$, there is a problem because k_{max} changes. Recall that the value of $k_{max} = \lceil \log_2 N \rceil + c$. In the border case if $N = 2^t - 1$ then $\lceil \log_2 N \rceil = t - 1$ and $\lceil \log_2(N+1) \rceil = t$, so we need to compare $\Psi(N, t-1+c)$ and $\Psi(N+1, t+c)$. In this case:

$$\begin{aligned} \frac{\Psi(N, t-1+c)}{\Psi(N+1, t+c)} &= \frac{\prod_{j=0}^{t+c-2} (1 - e^{-\frac{N}{2^{j+1}}}) \cdot \prod_{j=0}^{t-1} (1 - e^{-\frac{N+2}{2^{j+1}}})}{\prod_{j=0}^{t+c-2} (1 - e^{-\frac{N+1}{2^{j+1}}}) \cdot \prod_{j=0}^{t-1} (1 - e^{-\frac{N+1}{2^{j+1}}})} \\ &= \prod_{j=0}^{t+c-2} \frac{(1 - e^{-\frac{N}{2^{j+1}}}) (1 - e^{-\frac{N+2}{2^{j+1}}})}{(1 - e^{-\frac{N+1}{2^{j+1}}})^2} \cdot \left(\frac{1 - e^{-\frac{N+2}{2^{t+c}}} }{1 - e^{-\frac{N+1}{2^{t+c}}} } \right) \quad (4) \end{aligned}$$

While the product term in (4) is less than 1, the term in the big brackets is greater than 1. It is hard to decide whether the whole formula is less than 1 or not. Although we cannot compare $\Psi(2^t - 1, t - 1 + c)$ and $\Psi(2^t, t + c)$, in Lemma 6 we can show a weaker result (note 2^{t-1} in the lemma instead of $2^t - 1$):

Lemma 6. For all $t \in \mathbb{Z}^+$, $\Psi(2^{t-1}, t - 1 + c) < \Psi(2^t, t + c)$ where c is as defined in Lemma 4.

Lemma 6 is useful because combining it and Lemma 5, we can prove the following lemma:

Lemma 7. Let $t_0 \in \mathbb{Z}^+$, if $\Psi(2^{t_0}, t_0 + c) \geq e^{-\epsilon}$, then for any $N \geq 2^{t_0}$, $\epsilon > 0$, $\Psi(N, k_{max}) \geq e^{-\epsilon}$, where k_{max} is as defined in Lemma 4.

Now we are ready to state the next theorem:

Theorem 3. Let $\epsilon > 0$, and c, k_{max} as defined in Lemma 4. Let t_0 be the smallest positive integer that satisfies $\Psi(2^{t_0}, t_0 + c) \geq e^{-\epsilon}$. Let N_2 be the smallest integer in $(2^{t_0-1}, 2^{t_0}]$ such that $\Psi(N_2, t_0 + c) \geq e^{-\epsilon}$. Then, (1) $\forall N \geq N_2, k \leq k_{max}, \frac{p_{N,k}}{p_{N+1,k}} \geq e^{-\epsilon}$, and (2) $\forall N \geq N_2, k > k_{max}, p_{N,k} \leq \delta$ and $p_{N+1,k} \leq \delta$.

Computing N_0 Combining all the above together, we can use Algorithm 8 to compute N_0 for a given (ϵ, δ) pair:

Algorithm 8: Find $N_0(\epsilon, \delta)$

Input: $\epsilon > 0, 0 < \delta < 1$

Result: $N_0 \in \mathbb{Z}^+$

- 1 $k_{min} = \max(\lceil \log_2 \frac{1}{\epsilon} \rceil - 1, 0)$;
/* $1 - q_{N, k_{min}}$ decreases monotonically in N . */
 - 2 Starting from 1, use an exponential search in $[1, +\infty]$ to find N_1 that is the smallest integer satisfying
$$1 - q_{N_1, k_{min}} = 1 - \sum_{j=0}^{k_{min}} (-1)^{v(j)} e^{-\frac{j N_1}{2^{k_{min}}}} \leq \delta$$
;
 - 3 $c = \lceil \frac{-1 + \sqrt{1 + 8 \log_2 \frac{1}{\delta}}}{2} \rceil$;
 - 4 Starting from 1, use an exponential search in $[1, +\infty]$ to find t_0 that is the smallest integer satisfying
$$\prod_{j=0}^{k_{max}-1} \frac{(1 - e^{-\frac{2^j}{2^{j+1}}})}{(1 - e^{-\frac{2^j+1}{2^{j+1}}})} \geq e^{-\epsilon}$$
, where $k_{max} = t_0 + c$;
/* search backwardly in $(2^{t_0-1}, 2^{t_0}]$ */
 - 5 **for** $i = 2^{t_0}; i > 2^{t_0-1}; i --$ **do**
 - 6 **if** $(\prod_{j=0}^{t_0+c-1} \frac{(1 - e^{-\frac{N}{2^{j+1}}})}{(1 - e^{-\frac{N+1}{2^{j+1}}})} < e^{-\epsilon})$ **then**
 - 7 $N_2 = i + 1$;
 - 8 **break**;
 - 9 **end**
 - 10 **end**
 - 11 Output $N_0 = \max(N_1, N_2)$;
-

Regarding the algorithm, we have the following theorem:

Theorem 4. For all $\epsilon, \delta \in \mathbb{R}^+$ and $\delta \in (0, 1)$, let $N_0 = \text{Find}N_0(\epsilon, \delta)$. When all DPs use a single FM sketch, our protocol satisfies (ϵ, δ) -differential privacy if the cardinality of the union of all DP's set is greater or equal to N_0 .

The running time of Algorithm 8 is bounded by the search time, and in turn the values of, N_1 and N_2 . We have the following Theorem:

Theorem 5. In algorithm 8, N_1 and N_2 increase monotonically as the parameter ϵ or δ decrease.

Therefore for smaller (ϵ, δ) , the algorithm will take longer to run. However this will not be a problem in practice. As an example, we ran Algorithm 8 with extremely small parameters $\epsilon = 2^{-40}$ and $\delta = 2^{-80}$, $N_0 = \max(N_1, N_2)$ found by the algorithm is 30,865,997,083,798, and the running time was in the order of seconds⁴. Therefore, for all (ϵ, δ) normally used in practice, N_1, N_2 will not be too large and the algorithm can be efficiently computed (see also Table 1, in which the values were computed with $(\frac{\epsilon}{m}, \frac{\delta}{m})$).

5.2.3 Find N_0 : Multiple Sketches

The Bound By Composition Theorems If the DPs use m FM sketches, then the output of the protocols is $Z_N = z_{N,1} + \dots + z_{N,m}$, where $z_{N,i}$ is extracted from the i -th FM sketch. The input set encoded by each FM sketch is the same, i.e. the union of observations from all DPs, and the hash keys are different. Therefore, using m FM sketches is like querying a privacy mechanism m times, and the randomization of the mechanism is independent for each query. The basic composition theorem (Theorem 3.16, [26]) states that if the base differential privacy mechanism is (ϵ_0, δ_0) -differentially private, then after m queries, any function of the m query results is at least $(m\epsilon_0, m\delta_0)$ -differentially private. Therefore, in the m sketches case, given the target (ϵ, δ) we want to achieve, it suffices if each single FM sketch satisfies $(\frac{\epsilon}{m}, \frac{\delta}{m})$ -differential privacy. When m is large, the advanced composition theorem (Theorem 3.20, [26]) gives a better bound. For a base mechanism that is (ϵ_0, δ_0) differentially private, after m queries, the result is at least $(\epsilon, m\delta_0 + \delta')$ -differential privacy, where

$$\epsilon = \sqrt{2m \ln \frac{1}{\delta'}} \epsilon_0 + m\epsilon_0(e^{\epsilon_0} - 1), \text{ for any } \delta' > 0. \quad (5)$$

Hence, given the target (ϵ, δ) , we can obtain (ϵ_0, δ_0) , then an initial bound $\widehat{N}_0 = \text{find}N_0(\epsilon_0, \delta_0)$. For all $N \geq \widehat{N}_0$, (ϵ, δ) -differential privacy holds, due to Theorem 4 and the composition theorems.

In Table 1, we show some \widehat{N}_0 for different combinations of parameters. When $m = 100$, the basic composition theorem gives better results, so we set $\epsilon_0 = \frac{\epsilon}{100}$, $\delta_0 = \frac{\delta}{100}$. For all other m , we obtain ϵ_0, δ_0 through the advanced composition theorem. We simply set $\delta' = \frac{\delta}{2}$ and $\delta_0 = \frac{\delta}{2m}$, then we can get ϵ_0 by (5). Note that in the table, when $m = 100$, we get the same \widehat{N}_0 in the cases when $\epsilon = 0.2$ and $\epsilon = 0.3$. This is because in both cases $N_1 > N_2$, so $\widehat{N}_0 = N_1$. The value of N_1 is a function

⁴The implementation is based on Arb (<http://arblib.org/>), a C library supporting arbitrary precision real arithmetic.

$\epsilon \backslash m$	100	1000	2000	4000
1	2053	4596	9387	9564
0.5	4123	9210	18791	19146
0.3	8261	18437	37601	38310
0.2	8261	36891	37601	76638
0.1	16538	73800	75219	153295

Table 1: The value of \widehat{N}_0 for different ϵ , m and fixed $\delta = 2^{-40}$, $w = 32$

of $\lceil \log_2 \frac{1}{\epsilon_0} \rceil - 1$ and δ_0 . The same δ_0 is used in both cases and $\lceil \log_2 \frac{100}{0.2} \rceil - 1 = \lceil \log_2 \frac{100}{0.3} \rceil - 1$, so the algorithm gives the same \widehat{N}_0 . For the same reason, we get the same \widehat{N}_0 for $\epsilon = 0.2$ and $\epsilon = 0.3$ when $m = 2000$.

The bound \widehat{N}_0 by composition theorems is rather loose and can be further improved. Next we will first show how to compute the PMF of Z_N , then how we can get an improved bound N_0 computationally.

PMF: m FM sketches The PMF of Z_N can be obtained through the probability generating functions (pgf for short) [33]. We know that the pgf of a discrete random variable X taking values in non-negative integer $[0, j]$ is defined as:

$$G_X(t) = \mathbf{E}(t^X) = \sum_{k=0}^j Pr[X = k] \cdot t^k.$$

Therefore for $z_{N,i}$, the pgfs are:

$$G_{z_{N,i}}(t) = \sum_{k=0}^{w-1} p_{N,k} \cdot t^k.$$

We use pgfs here because they are particularly useful for dealing with the sum of independent random variables. In fact, for $Z_N = \sum_{i=1}^m z_{N,i}$, the pgf is:

$$G_{Z_N}(t) = (G_{z_{N,i}}(t))^m = \left(\sum_{k=0}^{w-1} p_{N,k} \cdot t^k \right)^m. \quad (6)$$

Another property of a pgf is that the PMF of X can be recovered by taking derivatives of $G_X(t)$:

$$Pr[X = k] = \frac{G_X^{(k)}(0)}{k!}. \quad (7)$$

Expanding $G_{Z_N}(t)$, we will get the $m(w-1)$ -th degree polynomials $\sum_{K=0}^{m(w-1)} a_K t^K$, where a_K are coefficients and t is the indeterminate. Then by (7), we have:

$$Pr[Z_N = K] = \frac{G_{Z_N}^{(K)}(0)}{K!} = a_K. \quad (8)$$

Refining the Bound In the m FM sketches case, (ϵ, δ) -differential privacy holds if for every $0 \leq K \leq m(w-1)$:

$$e^{-\epsilon} \cdot Pr[Z_{N+1} = K] - \frac{\delta}{e^\epsilon} \leq Pr[Z_N = K] \leq e^\epsilon \cdot Pr[Z_{N+1} = K] + \delta. \quad (9)$$

Therefore, we can use algorithm 9 to find the improved N_0 .

Algorithm 9: *RefineBound*($\epsilon, \delta, \widehat{N}_0, m, w$)

Input: $\epsilon, \delta \in \mathbb{R}^+$ and $\delta \in (0, 1), \widehat{N}_0, m, w \in \mathbb{Z}^+$
Result: $N_0 \in \mathbb{Z}^+$

```

1 stop = false;
2  $N_0 = \widehat{N}_0 + 1$ ;
3 do
4    $N_0 = N_0 - 1$ ;
5   Compute the polynomials  $G_{Z_{N_0}}(t)$  and  $G_{Z_{N_0-1}}(t)$  using (6);
6   for  $K = 0; K \leq m(w-1); K++$  do
7     Let  $Pr[Z_N = K]$  be the coefficient of the  $K$ -th degree term
      of  $G_{Z_{N_0-1}}(t)$ ;
8     Let  $Pr[Z_{N+1} = K]$  be the coefficient of the  $K$ -th degree
      term of  $G_{Z_{N_0}}(t)$ ;
9     if  $Pr[Z_N = K]$  and  $Pr[Z_{N+1} = K]$  don't satisfy (9) then
10      stop = true;
11      break;
12   end
13 end
14 while stop = false and  $N_0 > 0$ ;
15 output  $N_0$ 

```

Algorithm 9 starts from \widehat{N}_0 and computationally verifies $N < N_0$ backwardly. It stops at N_0 when $N_0 - 1$ does not satisfy differential privacy anymore. This N_0 is the improved bound and it is guaranteed that for all $N \geq N_0$, our protocol satisfies (ϵ, δ) -differential privacy at the given (m, w) parameters. In Table 2, we show the improved bound computed from Algorithm 9. Compared to the values in Table 1, the improved bound is significantly better.

The running time of Algorithm 9 is dominated by Step 5, in which the pgfs are computed. Computing pgfs involving polynomial exponentiation and the time increases when m increases. For example, to get numbers in Table 2, it took 78 ms, 5350 ms, 21468 ms and 90237 ms to compute a single $G_{Z_{N_0}}(t)$ when $m = 100, 1000, 2000, 4000$ respectively. When \widehat{N}_0 is large, backward verification by Algorithm 9 could take quite long time. That said, it should be noted that this verification needs only to be done once for each parameter combination.

$\epsilon \backslash m$	100	1000	2000	4000
1	85	254	355	497
0.5	166	496	693	969
0.3	273	813	1136	1587
0.2	404	1205	1682	2351
0.1	790	2359	3293	4600

Table 2: The value of N_0 by Algorithm 9 for different ϵ , m and fixed $\delta = 2^{-40}$, $w = 32$

The bound N_0 can easily be achieved in real world applications. For example, when $\epsilon = 0.3$, which is recommended for safe measurements in anonymity networks [38], even with a large $m = 4000$, N_0 is only 1587. For smaller ϵ values, N_0 are still reasonably small across different m values. Note that N_0 is the lower bound, therefore the privacy level is guaranteed even if the actual cardinality is larger than N_0 . We can also see that for the same privacy parameters, a larger set allows us to get a better accuracy (by allowing a larger m at the same privacy level). This means we can get both good utility and good privacy if the set is large.

6 Experimental Evaluation

We have implemented a prototype of our protocol in C++. The source code of the protocol is available online⁵. We used the implementation of Overdrive (low gear) in the SPDZ2 repository⁶ for the pre-processing part, and implemented our offline and online protocols on top of that. We compare the performance of our protocol to the state of the art [30]. The implementation of [30] provided by the authors is in Go and does not fully support multi-threading. For a fair comparison, we re-implemented the protocol in [30] in C++. In this implementation, we use OpenSSL 1.0.1 for all cryptographic operations and pthread for multi-threading. The performance of our new implementation is much better than that reported in [30]. We used 40 for the statistical security parameter and 128 for the computational security parameter in all experiments.

We ran all CPs in Amazon AWS. We used the EC2 instance type r5.4xlarge (on-demand) for each CP. Each instance has 16 vCPUs (8 physical cores) based on Intel Xeon Platinum 8000 series (Skylake-SP) CPUs, 128GB RAM, one network interface up to 10 Gbps LAN speed, and costs \$1.008 - \$1.12 per hour in US data centers. We conducted experiments both in a LAN environment (all CPs were in the Oregon AWS data center), and a WAN environment (CPs were distributed in 4 different AWS data centers in the US⁷). The DPs ran on desktops, with a typical hardware configuration of an Intel Quadcore i7-6700k CPU and 16 GB RAM. We used 20 DPs in all experiments and varied the number of CPs.

In Table 3, we show the total running time and communication (send+receive) cost of our protocol in the offline and online phases. We implemented the group authenticated key exchange protocol in [39]. The offline phase measurement includes the costs of the SPDZ pre-processing protocol and our offline protocols. The online phase measurement includes all online protocols, from DP sharing the sketch to the CPs outputting Z_N (using the ExtractZ protocol in LAN and ExtractZBS protocol in WAN). Note we do not include the time used by the DPs to collect data because this time is irrelevant

to our protocol. In the experiments, the DPs first did the initial sharing and then immediately the final sharing of the Oblivious FM sketches. The running time and communication cost shown in the table are the average of those measured over all CPs. For the running time, we show the time measured in LAN and WAN. The communication costs in LAN and WAN are almost the same, thus we only show the larger one of the two. We varied the number of distinct elements in the experiments, from 20000, to 1 million (10^6), to 1 billion (10^9). This change affects the size of the modulus p (55, 60, 70 respectively) and the size of the sketches w (19, 24, 34 respectively). We also used different number of sketches (m) for different accuracy levels. As we can see in the table, the total running time is dominated by the offline phase. While the offline running time is in the order of minutes, the online running time is only in the order of seconds. We can also see that the offline running time is less than 1 hour even with the largest parameter group, and since the offline computation can be done during the period when the DPs are collecting the data, the performance should be acceptable (many applications may only require daily or even less frequent update of the estimate). The protocol has good scalability: when N increases from 20000 to 10^9 (50000 times), the running time increases only to about 2 times ($\log(10^9)/\log(20000) \approx 2$). The running time in LAN is much less than that measured in WAN. The differences in network bandwidth and latency are likely the causes of the slowdown. Communication-wise, the offline phase cost is much higher than the online phase cost. As we can see in Table 4, most of the cost in the offline phase is due to the Triple generation protocol in SPDZ, which utilizes heavy machinery such as somewhat homomorphic encryption and zero-knowledge proofs. In Table 4, we also show the differences in performance for the ExtractZ (Protocol 7) and ExtractZBS (Protocol 10, Appendix A). The results confirm that in the high network latency setting, ExtractZBS performs better due to fewer communication rounds/interactions.

As a comparison, we show in Figure 4 the total running time and communication cost of the protocol in [30]. In the experiments, we used 5 CPs (16 threads) and 20 DPs. We varied N from 20000 to 50000, and as in [30], set the number of bins to $10 \cdot N$ so the collision probability is less than 10%. We also set (ϵ, δ) for differential privacy to $(0.3, 10^{-12})$, the default values used in [30]. Note the parameters are weaker than those for our protocol: with $N = 20000$ and other parameters in the experiments, our protocol can easily achieve $(0.1, 10^{-12})$ -differential privacy, and even better privacy when N grows bigger. We only tested with all CPs in the same LAN, as the figures in the WAN setting would be even higher. As we can see, the protocol in [30] is much slower than ours, and its running time increases much faster. When $N = 20000$, its running time in LAN is about 1.2 times of ours in WAN, and 8.5 times of ours in LAN (both $m = 4000$); when $N = 50000$, it needs almost 2.5 hours in LAN, while our protocol (in WAN) with $N = 10^9$ only needs less than 50 minutes (offline+online).

⁵<https://github.com/saftoes/pdce>

⁶<https://github.com/bristolcrypto/SPDZ-2>

⁷N Virginia, Ohio, Northern California, Oregon.

			$N = 20000$			$N = 10^6$			$N = 10^9$		
			m=1000	m=2000	m=4000	m=1000	m=2000	m=4000	m=1000	m=2000	m=4000
Running Time (s)	LAN	Offline	66.5	132.2	222.8	78.1	154.7	307.6	149.3	257.3	515.8
		Online	0.079	0.151	1.997	0.110	0.189	0.271	0.201	0.377	0.522
	WAN	Offline	320	624.1	1470.5	411.7	811.1	1578.9	757.1	1421.8	2944.2
Online		2.414	2.036	2.623	1.754	2.360	2.934	2.689	3.031	5.026	
Communication (GB)	Offline		10.7	21.4	35.2	12.09	24.18	48.5	23.3	39.05	78.3
	Online		0.008	0.016	0.031	0.010	0.020	0.041	0.028	0.056	0.120

Table 3: Total running time and communication cost: 5 CPs (16 threads), 20 DPs.

		Running Time (s)		Comm. (GB)
		LAN	WAN	
Group AKE (per DP)		0.014	0.46	6.36×10^{-6}
Offline	Triple	417.0	2414.1	68.5
	Rand	50.6	452.4	7.4
	Rand2	47.4	70.3	1.83
	RandExp	0.8	7.4	0.61
Online	Share (per DP)	0.155	1.877	0.0087
	MergeShare	0.00130	0.00129	N/A
	ZeroTest	0.32	2.345	0.070
	ExtractZ	0.049	1.482	0.034
	ExtractZBS	0.063	0.803	0.042

Table 4: Performance breakdown: 5 CPs (16 threads), 20 DPs, $N = 10^9$, $m = 4000$

The running time of [30] is slightly convex due to a quadratic step in a zero-knowledge proof sub-protocol. The communication complexity of the protocol in [30] is linear. When N is small, the protocol in [30] has a much smaller communication cost compared to ours, e.g. 1.4 GB vs 35.2 GB when $N = 20000$. However since the communication complexity of the protocol in [30] is linear and that of ours is logarithmic, the communication cost of the protocol in [30] will exceed that of ours eventually. As an estimation, when N is 10^6 , the communication cost of the protocol in [30] would be 60 GB roughly, which is already higher than ours (48.5 GB).

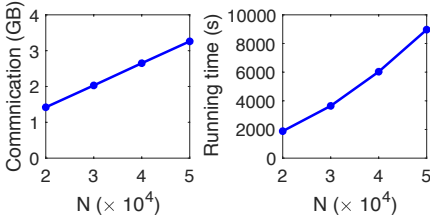
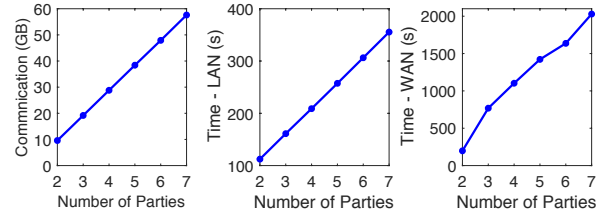


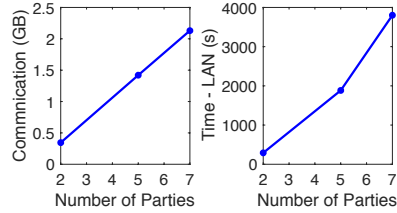
Figure 4: Performance of protocol in [30] (LAN)

Next, we show in Figure 5 the performance of our protocol and the protocol in [30] with a varying number of CPs. For our protocol, we fixed N to 10^9 and m to 4000, with a varying number of CPs from 2 to 7. As we can see, the communication cost and the running time in LAN increase linearly in the number of CPs. The line of the running time in WAN is not very regular, but we can see that the running time is roughly linear. In typical applications, the number of CPs is quite unlikely to exceed 10. However in the case of more CPs, we could switch the SPDZ pre-processing protocol to High Gear. High Gear’s performance surpasses Low Gear (we currently

use) when executed with a high number of parties (more than 10 as reported in [41]). As the computation time of our protocol is dominated by the SPDZ pre-processing, this would allow us to handle more CPs more gracefully. For the protocol in [30], we fixed N to 20000, and used 2, 5, 7 CPs in the experiment. The communication cost of this protocol is also linear in the number of CPs, but the running time is slightly worse than linear. The results are consistent with those reported in [30].



(a) Our protocol, $N = 10^9$, $m = 4000$



(b) Protocol in [33], $N = 20000$

Figure 5: Performance with different number of CPs

In Figure 6, we show the distribution of the relative errors $\left(\frac{|\tilde{N}-N|}{N}\right)$ where \tilde{N} is the cardinality estimated from the sketches and N is the true cardinality) when using a different number of FM sketches. We used $m = 1000, 2000$ and 4000 sketches, using two sets with 20000 and 10^6 random elements as inputs. We repeated each experiment 1000 times and drew the histograms. As we can see, when m increases, the max relative error decreases, and the distribution gets more concentrated towards 0. With $m = 4000$, about 99% of the estimations have a relative error less than 3%, and the maximum relative error observed was 4.3%. On the other hand, the estimations using the method in [30] had a slightly higher relative error (see the full version) due to the hash collisions and the noise added to achieve differential privacy.

Since the cardinality count produced by the protocol in [30] is also approximate, it would be interesting to see whether our differential privacy analysis can result in a cheaper variant of that protocol, and if so how would the performance of

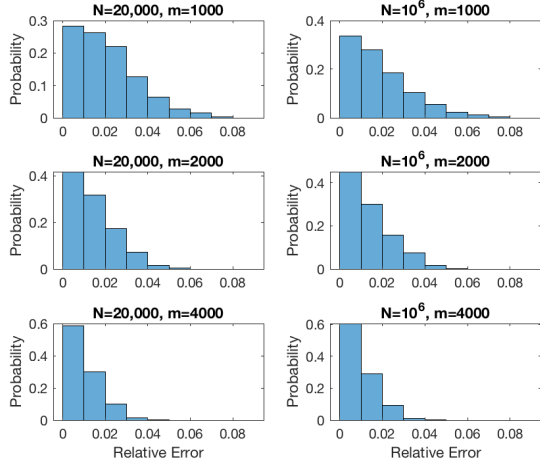


Figure 6: Distribution of relative errors the variant compare to that of our protocol. In principle the protocol of [30] could also obtain differential privacy for free with honest DPs and a private hash key, although we have not done the analysis and the analysis may not be trivial. If [30] achieves differential privacy by hashing, then the CPs do not need to add noise. However, the performance improvement would be around 20-30% at most, based on our experience of implementing the protocol. The performance would be in the same order as it is now, and thus still much worse than that of our protocol. This is because the main factors affecting the performance of [30] are not adding noise but (1) public key encryption; (2) verifiable shuffling and zero-knowledge proofs; (3) superlinear (in the maximum measurable cardinality) computational and communication complexity.

7 Conclusion and Future Work

In this paper, we present and analyse a PDCE protocol. The protocol is efficient and scalable, due to the use of FM sketches as the underlying data structure for cardinality estimation, and the use of efficient secret sharing based MPC primitives. We proved the security of the protocol against a malicious adversary in the UC framework. More interestingly, we showed that the combination of secure computation and the FM sketches allows us to get (ϵ, δ) -differential privacy for free. We implemented our protocol and evaluated it experimentally. Our experiments showed that the protocol is much more efficient and scalable than the state of the art [30].

We would like to continue investigating the use of data structures in secure computation protocols to improve their efficiency and scalability. Data structures such as sketches could lead to sub-linear complexity protocols, which are highly desirable for Big Data applications. We would also like to investigate the relationship between differential privacy and sketches, to extend and generalize the results in this paper to other sketches/data structures.

Acknowledgement

We thank shepherd Mathias Lécuyer as well as the anonymous reviewers for their insightful comments. This research was supported in part by UK EPSRC under grant EP/M013561/2; National Natural Science Foundation of China under grant 61722203 (Outstanding Youth Foundation), U1936218 (Joint Fund Project), 62032012, 62072132, and 61771259; National Key Research and Development Program of China under grant 2020YFB1005700.

References

- [1] Health Insurance Portability and Accountability Act of 1996. <https://aspe.hhs.gov/report/health-insurance-portability-and-accountability-act-1996>, 1996.
- [2] Gramm-Leach-Bliley Act. <https://www.ftc.gov/tips-advice/business-center/privacy-and-security/gramm-leach-bliley-act>, 1999.
- [3] General Data Protection Regulation. <https://eur-lex.europa.eu/legal-content/EN/TXT/HTML/?uri=CELEX:32016R0679&from=EN>, 2018.
- [4] The Royal Society Report on Privacy Enhancing Technologies. <https://royalsociety.org/-/media/policy/projects/privacy-enhancing-technologies/privacy-enhancing-technologies-report.pdf>, 2019.
- [5] Gergely Ács and Claude Castelluccia. A case study: privacy preserving release of spatio-temporal density in paris. In *KDD*, pages 1679–1688, 2014.
- [6] Sameer Agarwal, Barzan Mozafari, Aurojit Panda, Henry Milner, Samuel Madden, and Ion Stoica. Blinkdb: queries with bounded errors and bounded response times on very large data. In *EuroSys*, pages 29–42, 2013.
- [7] Akamai. Real-Time Web Metrics Methodology. <https://www.akamai.com/uk/en/resources/visualizing-akamai/real-time-web-monitor/real-time-web-metrics-methodology.jsp>.
- [8] Vikas G. Ashok and Ravi Mulkamala. A scalable and efficient privacy preserving global itemset support approximation using bloom filters. In *DBSec*, pages 382–389, 2014.
- [9] Martin Azizyan, Ionut Constandache, and Romit Roy Choudhury. Surroundsense: mobile phone localization via ambience fingerprinting. In *MOBICOM*, pages 261–272, 2009.
- [10] Donald Beaver. Efficient multiparty protocols using circuit randomization. In *CRYPTO*, pages 420–432, 1991.
- [11] Zvika Brakerski, Craig Gentry, and Vinod Vaikuntanathan. (leveled) fully homomorphic encryption without bootstrapping. In *ITCS*.
- [12] Ran Canetti. Universally composable security: A new paradigm for cryptographic protocols. In *Proceedings 2001 IEEE International Conference on Cluster Computing*, pages 136–145, 2001.
- [13] Ran Canetti. Universally composable security: A new paradigm for cryptographic protocols. In *FOCS*, pages 136–145, 2001.

- [14] T.-H. Hubert Chan, Mingfei Li, Elaine Shi, and Wenchang Xu. Differentially private continual monitoring of heavy hitters from distributed streams. In *PETS*, pages 140–159, 2012.
- [15] Graham Cormode. Data sketching. *Commun. ACM*, 60(9):48–55, August 2017.
- [16] Ronald Cramer, Ivan Damgård, and Jesper Buus Nielsen. *Secure Multiparty Computation and Secret Sharing*. Cambridge University Press, 2015.
- [17] Emiliano De Cristofaro, Paolo Gasti, and Gene Tsudik. Fast and private computation of cardinality of set intersection and union. In *CANS*, pages 218–231, 2012.
- [18] Ivan Damgård, Marcel Keller, Enrique Larraia, Valerio Pastro, Peter Scholl, and Nigel P. Smart. Practical covertly secure MPC for dishonest majority - or: Breaking the SPDZ limits. In *ESORICS*, pages 1–18, 2013.
- [19] Ivan Damgård, Valerio Pastro, Nigel P. Smart, and Sarah Zarkarias. Multiparty computation from somewhat homomorphic encryption. In *CRYPTO*, pages 643–662, 2012.
- [20] Alex Davidson and Carlos Cid. An efficient toolkit for computing private set operations. In *ACISP*, pages 261–278, 2017.
- [21] Damien Desfontaines, Andreas Lochbihler, and David A. Basin. Cardinality estimators do not preserve privacy. In *PETS*.
- [22] Roger Dingledine, Nick Mathewson, and Paul F. Syverson. Tor: The second-generation onion router. In *USENIX Security*, pages 303–320, 2004.
- [23] Changyu Dong and Grigorios Loukides. Approximating private set union/intersection cardinality with logarithmic complexity. *IEEE Trans. Information Forensics and Security*, 12(11):2792–2806, 2017.
- [24] Cynthia Dwork. Differential privacy. In *ICALP*, pages 1–12, 2006.
- [25] Cynthia Dwork, Krishnaram Kenthapadi, Frank McSherry, Ilya Mironov, and Moni Naor. Our data, ourselves: Privacy via distributed noise generation. In *EUROCRYPT*, pages 486–503, 2006.
- [26] Cynthia Dwork and Aaron Roth. The algorithmic foundations of differential privacy. *Found. Trends Theor. Comput. Sci.*, 9(3-4):211–407, 2014.
- [27] Rolf Egert, Marc Fischlin, David Gens, Sven Jacob, Matthias Senker, and Jörn Tillmanns. Privately computing set-union and set-intersection cardinality via bloom filters. In *ACISP*, pages 413–430, 2015.
- [28] Tariq Elahi, George Danezis, and Ian Goldberg. Privex: Private collection of traffic statistics for anonymous communication networks. In *ACM CCS*, pages 1068–1079, 2014.
- [29] Úlfar Erlingsson, Vasily Pihur, and Aleksandra Korolova. RAP-POR: randomized aggregatable privacy-preserving ordinal response. In *ACM CCS*, pages 1054–1067, 2014.
- [30] Ellis Fenske, Akshaya Mani, Aaron Johnson, and Micah Sherr. Distributed measurement with private set-union cardinality. In *ACM CCS*, pages 2295–2312, 2017.
- [31] Philippe Flajolet and G. Nigel Martin. Probabilistic counting algorithms for data base applications. *J. Comput. Syst. Sci.*, 31(2):182–209, 1985.
- [32] Oded Goldreich. *The Foundations of Cryptography*. Cambridge University Press, 2004.
- [33] Geoffrey Grimmett and David Stirzaker. *Probability and random processes*. Oxford University Press, third edition edition, 2001.
- [34] Hazar Harmouch and Felix Naumann. Cardinality estimation: An experimental survey. *PVLDB*, 11(4):499–512, 2017.
- [35] Rob Harrison, Qizhe Cai, Arpit Gupta, and Jennifer Rexford. Network-wide heavy hitter detection with commodity switches. In *SOSR*, 2018.
- [36] Carmit Hazay, Yuval Ishai, and Muthuramakrishnan Venkatasubramanian. Actively secure garbled circuits with constant communication overhead in the plain model. In *TCC*, pages 3–39, 2017.
- [37] Stefan Heule, Marc Nunkesser, and Alexander Hall. Hyper-loglog in practice: algorithmic engineering of a state of the art cardinality estimation algorithm. In *EDBT*, pages 683–692, 2013.
- [38] Rob Jansen and Aaron Johnson. Safely measuring tor. In *ACM CCS*, pages 1553–1567, 2016.
- [39] Jonathan Katz and Moti Yung. Scalable protocols for authenticated group key exchange. In *CRYPTO*, pages 110–125, 2003.
- [40] Marcel Keller, Emmanuela Orsini, and Peter Scholl. MASCOT: faster malicious arithmetic secure computation with oblivious transfer. In *ACM CCS*, pages 830–842, 2016.
- [41] Marcel Keller, Valerio Pastro, and Dragos Rotaru. Overdrive: Making SPDZ great again. In *EUROCRYPT*, pages 158–189, 2018.
- [42] Helger Lipmaa and Tomas Toft. Secure equality and greater-than tests with sublinear online complexity. In *ICALP*, pages 645–656, 2013.
- [43] Akshaya Mani and Micah Sherr. Histore: Differentially private and robust statistics collection for tor. In *NDSS*, 2017.
- [44] Luca Melis, George Danezis, and Emiliano De Cristofaro. Efficient private statistics with succinct sketches. In *NDSS*, 2016.
- [45] Darakhshan J. Mir, S. Muthukrishnan, Aleksandar Nikolov, and Rebecca N. Wright. Pan-private algorithms via statistics on sketches. In *PODS*, pages 37–48, 2011.
- [46] David Moore, Geoffrey M. Voelker, and Stefan Savage. Inferring internet denial-of-service activity. In *USENIX Security*, 2001.
- [47] A. B. M. Musa and Jakob Eriksson. Tracking unmodified smartphones using wi-fi monitors. In *SenSys*, pages 281–294, 2012.
- [48] Nikos Ntarmos, Peter Triantafillou, and Gerhard Weikum. Counting at large: Efficient cardinality estimation in internet-scale data networks. In *ICDE*, 2006.
- [49] Information Commissioner’s Office. Wi-fi location analytics. 2016.
- [50] Chris Peikert, Vinod Vaikuntanathan, and Brent Waters. A framework for efficient and composable oblivious transfer. In *CRYPTO*, pages 554–571, 2008.
- [51] Vibhor Rastogi and Suman Nath. Differentially private aggregation of distributed time-series with transformation and encryption. In *SIGMOD*, pages 735–746, 2010.
- [52] Nathaniel Schenker and Trivellore E. Raghunathan. Combining information from multiple surveys to enhance estimation of measures of health. *Statistics in medicine*, 26(8):1802–1811, 2007.

- [53] Björn Scheuermann and Martin Mauve. Near-optimal compression of probabilistic counting sketches for networking applications. In *DIALM-POMC*, 2007.
- [54] Elaine Shi, T.-H. Hubert Chan, Eleanor G. Rieffel, Richard Chow, and Dawn Song. Privacy-preserving aggregation of time-series data. In *NDSS*, 2011.
- [55] Hagen Sparka, Florian Tschorsch, and Björn Scheuermann. P2KMV: A privacy-preserving counting sketch for efficient and accurate set intersection cardinality estimations. *IACR Cryptology ePrint Archive*, 2018:234, 2018.
- [56] Rade Stanojevic, Mohamed Nabeel, and Ting Yu. Distributed cardinality estimation of set operations with differential privacy. In *IEEE PAC*, pages 37–48, 2017.
- [57] STASTICA. Marks & Spencer average weekly footfall in the United Kingdom (UK) 2009-2018. <https://www.statista.com/statistics/413515/marks-and-spencer-mands-average-weekly-footfall-united-kingdom-uk/>.
- [58] Stephanie Clifford and Quentin Hardy. Attention, Shoppers: Store Is Tracking Your Cell. <https://www.nytimes.com/2013/07/15/business/attention-shopper-stores-are-tracking-your-cell.html>, 2013.
- [59] The Guardian. Shops can track you via your smartphone, privacy watchdog warns. <https://www.theguardian.com/technology/2016/jan/21/shops-track-smartphone-uk-privacy-watchdog-warns>, 2016.
- [60] Florian Tschorsch and Björn Scheuermann. An algorithm for privacy-preserving distributed user statistics. *Computer Networks*, 57(14):2775–2787, 2013.
- [61] Wei Xi, Jizhong Zhao, Xiang-Yang Li, Kun Zhao, Shaojie Tang, Xue Liu, and Zhiping Jiang. Electronic frog eye: Counting crowd using wifi. In *INFOCOM*, pages 361–369, 2014.
- [62] Qingjun Xiao, You Zhou, and Shigang Chen. Better with fewer bits: Improving the performance of cardinality estimation of large data streams. In *INFOCOM*, pages 1–9, 2017.

A Alternative Protocol for Extracting Estimator

Recall that z_N is the index of the first 0 bit in a sketch, thus extracting z_N can be converted to a search problem. Protocol 10 performs essentially a binary search. In Protocol 10, the bits in the sketch are first negated (lines 3 - 5). Then the sketch is divided into two halves. If all bits now in the first half are 0, then before negation, all of them were 1, which means the first 0 we are looking for is in the second half. Then we know z_N must be the size of the first half plus some offset into the second half, and we can throw away the first half and do a binary search on the second half to find the offset. If not all bits in the first half are 0, then the first 0 we are looking for is in the first half. Then we can throw away the second half and do another binary search on the first half. Obviously, we cannot reveal whether the first half is all 0 in the protocol, as this leaks information. So what we do is to sum all bits in the first half into x , then interpolate a lookup polynomial f such that $B_0 = f(x+1) = 1$ if $x = 0$ and

0 otherwise (lines 8 – 10). Then we obviously combine the first half and the second half, by multiplying every bit in the second half with B_0 and add the result to the first half (lines 12 – 17). Since the multiplications are independent, they can be batched together. Note also that an extra addition is needed if the two halves are not of the same size. If the first half is all 0, then we need to continue searching the second half. In this case, B_0 is 1 and what we get after the addition is the second half. If the first half is not all 0, then we do not have to search the second half at all. In this case B_0 is 0 and we get the first half after the addition. Then we start the while loop again until there are only few bits left to search. In this case, we take the bits left and do a lookup to finish the search (lines 20 – 22). There are $\log(w)$ iterations in the while loop, and in each iteration, we need two rounds: one round for line 10 (because of the multiplication in the *Lookup* protocol) and one round for the multiplications in the for loop starting at line 12.

Protocol 10: *ExtractZBS*($\llbracket \text{BFS}_{\cup}^1[0] \rrbracket, \dots, \llbracket \text{BFS}_{\cup}^1[w-1] \rrbracket, \dots, \llbracket \text{BFS}_{\cup}^m[0] \rrbracket, \dots, \llbracket \text{BFS}_{\cup}^m[w-1] \rrbracket$)

Input: $\llbracket \text{BFS}_{\cup}^1[0] \rrbracket, \dots, \llbracket \text{BFS}_{\cup}^1[w-1] \rrbracket, \dots, \llbracket \text{BFS}_{\cup}^m[0] \rrbracket, \dots, \llbracket \text{BFS}_{\cup}^m[w-1] \rrbracket$, the shares of the m binary FM sketches

Result: Z_N , the estimator extracted from the sketches

```

1  $\llbracket Z_N \rrbracket = 0$ ;
2 for  $i = 1; i \leq m; i++$  do
3   for  $j = 0; j \leq w-1; j++$  do
4      $\llbracket \text{BFS}_{\cup}^i[j] \rrbracket = 1 - \llbracket \text{BFS}_{\cup}^i[j] \rrbracket$ ; // negate the bit
5   end
6    $size = w, t = \lceil \frac{size}{2} \rceil, \llbracket z_{N,i} \rrbracket = 0$ ;
   // binary search until not worth it
7   while  $size > 3$  do
8      $\llbracket x+1 \rrbracket = 1 + \sum_{l=0}^{t-1} \llbracket \text{BFS}_{\cup}^i[l] \rrbracket$ ;
     // interpolate the lookup polynomial
9      $(t, \beta_0, \dots, \beta_t) \leftarrow \text{interpolate}()$ ;
     //  $B_0 = 1$  if  $x = 0$ ,  $B_0 = 0$  otherwise
10     $\llbracket B_0 \rrbracket = \text{Lookup}(\llbracket x+1 \rrbracket, t, \beta_0, \dots, \beta_t)$ ;
11     $\llbracket z_{N,i} \rrbracket = \llbracket z_{N,i} \rrbracket + t \cdot \llbracket B_0 \rrbracket$ ;
12    for  $j = 0; j < size - t; j++$  do
13       $\llbracket \text{BFS}_{\cup}^i[j] \rrbracket = \llbracket \text{BFS}_{\cup}^i[j] \rrbracket + \llbracket B_0 \rrbracket \cdot \llbracket \text{BFS}_{\cup}^i[j+t] \rrbracket$ 
      ;
14    end
15    if  $size$  is odd then
16       $\llbracket \text{BFS}_{\cup}^i[size-t] \rrbracket = \llbracket \text{BFS}_{\cup}^i[size-t] \rrbracket + \llbracket B_0 \rrbracket$ ;
17    end
18     $size = t, t = \lceil \frac{size}{2} \rceil$ ;
19  end
20   $\llbracket x \rrbracket = 1 + \sum_{i=0}^{size-1} 2^i \cdot \llbracket \text{BFS}_{\cup}^i[size-1] \rrbracket$ ;
  // interpolate the lookup polynomial
21   $(2^{size}, \beta_0, \dots, \beta_{2^{size}}) \leftarrow \text{interpolate}()$ ;
  // final lookup for the rest of the bits
22   $\llbracket z_{N,i} \rrbracket = \llbracket z_{N,i} \rrbracket + \text{Lookup}(\llbracket x \rrbracket, 2^{size}, \beta_0, \dots, \beta_{2^{size}})$ ;
23   $\llbracket Z_N \rrbracket = \llbracket Z_N \rrbracket + \llbracket z_{N,i} \rrbracket$ ;
24 end
25 return  $Z_N \leftarrow \text{Output}(\llbracket Z_N \rrbracket)$ ;

```
