
Differentially Private Linear Sketches: Efficient Implementations and Applications

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Linear sketches have been widely adopted to process fast data streams, and they
2 can be used to accurately answer frequency estimation, approximate top K items,
3 and summarize data distributions. When data are sensitive, it is desirable to
4 provide privacy guarantees for linear sketches to preserve private information while
5 delivering useful results with theoretical bounds. We show that linear sketches
6 can ensure privacy and maintain their unique properties with a small amount of
7 noise added at initialization. From the differentially private linear sketches, we
8 showcase that the state-of-the-art quantile sketch in the turnstile model can also be
9 private and maintain high performance. Experiments further demonstrate that our
10 proposed differentially private sketches are quantitatively and qualitatively similar
11 to noise-free sketches with high utilization on synthetic and real datasets.

12 1 Introduction

13 Data sketches are fundamental tools for data analysis, statistics, and machine learning [Cormode and
14 Yi, 2020]. Two of the most widely studied problems in data summaries are frequency estimation and
15 quantile approximation. Many real world applications need to estimate the frequency of each item
16 in the database and understand the overall distribution of the database. These applications include
17 stream processing [Das et al., 2009, Bailis et al., 2017], database management [Misra and Gries, 1982,
18 Metwally et al., 2005, Zhao et al., 2021a], caching [Zakhary et al., 2020], network monitoring [Gupta
19 et al., 2016, Ivkin et al., 2019], federated learning [Rothchild et al., 2020], among others.

20 On one hand, the motivation for data sketch algorithms is to efficiently process a large database
21 and extract useful knowledge, since computing the exact information for a large amount of data is
22 both time and memory intensive. For instance, Munro and Paterson [1980] proved that to find the
23 true median of a database with n items using p sequential passes requires at least $\Omega(n^{1/p})$ memory.
24 On the other hand, to protect user-level privacy, privacy-preserving algorithms limit the disclosure
25 of private information in the database so that an observer cannot infer much about an individual.
26 Recent works have shown that data sketches can be integrated with privacy-enhancing technologies
27 to provide insightful information and preserve individual privacy at the same time [Cormode, 2022].

28 Differential privacy [Dwork et al., 2006] is a widely-accepted definition of privacy. Recently, re-
29 searchers have observed that some data sketches are inherently differentially private [Blocki et al.,
30 2012, Smith et al., 2020], while many other data sketches need modifications to the algorithm to be
31 differentially private. In particular, a substantial amount of literature has focused on differentially
32 private data sketches for tasks such as linear algebra [Upadhyay, 2014, Arora et al., 2018], cardi-

33 nality estimation [Mir et al., 2011, Pagh and Stausholm, 2020, Dickens et al., 2022] and quantile
34 approximation [Tzamos et al., 2020, Gillenwater et al., 2021, Alabi et al., 2022].

35 In this paper, we introduce new differentially private algorithms that support both insertions and
36 deletions for frequency and top k estimation and quantile approximation. While many data sketches
37 assume an insertion-only model [Greenwald and Khanna, 2001, Shrivastava et al., 2004, Karnin
38 et al., 2016] or a bounded-deletion model [Jayaram and Woodruff, 2018, Zhao et al., 2021a,b], our
39 algorithms build on top of linear sketches [Charikar et al., 2002, Cormode and Muthukrishnan, 2005]
40 and operate in the most general turnstile model, which allows an arbitrary number of insertions
41 and deletions into the database. Earlier, researchers attempted to prove Count-Median [Charikar
42 et al., 2002] itself preserves differential privacy, but the authors acknowledged that there are issues
43 in the proof [Li et al., 2019]. Instead of proving that linear sketches, i.e., both Count-Median and
44 Count-Min, are inherently differentially private, we add a small amount of Gaussian noise at their
45 initialization to provide a privacy guarantee while maintaining linear sketches’ original properties,
46 providing high utility for frequency and top K estimations, and keeping update and query algorithms
47 unchanged. In addition, we propose the first differentially private quantile sketch in the turnstile
48 model by leveraging the differentially private linear sketch. Our differentially private sketches can
49 be queried an arbitrary number of times without affecting privacy guarantees based on the post-
50 processing immunity. Moreover, in our work, we removed the assumption on the existence of secret
51 uniform random hash functions to provide randomness in the algorithm. Following prior works [Choi
52 et al., 2020, Smith et al., 2020], we assume ideal random hash functions exist, and this assumption
53 can be replaced in practice by cryptographic hash functions [Dickens et al., 2022].

54 2 Preliminaries

55 Consider a database $X = \{i_t\}_{t \in [N]}$ of N items that are drawn from a large *universe* of size U , such
56 as IPv4 address of size 2^{32} , and for each insert or delete operations, one item can be inserted into
57 or deleted from the database X . To support ordered statistic such as quantile, we assume that the
58 *universe* is some finite totally ordered data universe.

59 **Definition 2.1.** Given a database X , the frequency of an item x is $f(x) = \sum_{t=1}^N \pi(i_t = x)$ where π
60 returns 1 if i_t is x , and 0 otherwise.

61 **Definition 2.2.** Given a database X of items drawn from an ordered universe, the rank of an item x
62 is $R(x) = \sum_{t=1}^N \pi(i_t \leq x)$ where π returns 1 if i_t is less or equal to x and 0 otherwise.

63 Given the large size of N , calculating the actual statistics, such as frequency and quantile, is often
64 hard, and hence most applications are satisfied with an *approximation*. The *randomized frequency*
65 *estimation problem* takes an accuracy parameter γ and a failure probability β such that, for any item
66 x , $|\hat{f}(x) - f(x)| \leq \gamma \cdot N$ with high probability $1 - \beta$, where $\hat{f}(x)$ is the estimated frequency and $f(x)$
67 is the true frequency [Cormode and Hadjieleftheriou, 2008]. In addition, the *randomized quantile*
68 *approximation problem* also takes an accuracy parameter γ and a failure probability β such that, for
69 any item x , $|\hat{R}(x) - R(x)| \leq \gamma \cdot N$ with high probability $1 - \beta$ where $\hat{R}(x)$ is the estimated rank
70 and $R(x)$ is the actual rank [Karnin et al., 2016].

71 2.1 Differential Privacy

72 **Definition 2.3.** Databases X and X' are neighbors ($X \sim X'$), if they differ in at most one element.

73 Through this paper, we assume the *update/replace* definition of differential privacy instead of
74 *add/remove* definition of differential privacy, in which one item in X is updated or replaced by
75 another item in X' [Vadhan, 2017].

76 **Definition 2.4** (Differential Privacy [Dwork et al., 2006]). A randomized algorithm M satisfies
77 (ϵ, δ) -differential privacy $((\epsilon, \delta)$ -DP) if for all neighboring databases X, X' and for all possible
78 events E in the output range of M , we have

$$\mathbb{P}(M(X) \in E) \leq e^\epsilon \cdot \mathbb{P}(M(X') \in E) + \delta.$$

79 When $\delta = 0$, ϵ -DP is known as pure DP, and when $\delta > 0$, (ϵ, δ) -DP is known as approximate DP.

80 **Definition 2.5** (Gaussian Mechanism [Dwork et al., 2006]). *Define the ℓ_2 sensitivity of a function*
 81 *$f : \mathbb{N}^{\mathcal{X}} \mapsto \mathbb{R}^d$ as*

$$\Delta_2(f) = \sup_{\text{neighboring } X, X'} \|f(X) - f(X')\|_2.$$

82 *The Gaussian mechanism \mathcal{M} with noise level σ is then given by*

$$\mathcal{M}(X) = f(X) + \mathcal{N}(0, \sigma^2 I_d).$$

83 Specifically, the Gaussian mechanism is known to satisfy a stronger notion of privacy known as zero-
 84 concentrated differential privacy (zCDP, defined below); zCDP lies between pure and approximate
 85 DP and can be translated into standard DP notations, as shown in Lemma 2.9. Moreover, zCDP
 86 satisfies cleaner composition theorems, as shown in Lemma 2.7.

Definition 2.6 (zCDP [Dwork and Rothblum, 2016, Bun and Steinke, 2016]). *A randomized mechanism M satisfies ρ -Zero-Concentrated Differential Privacy (ρ -zCDP), if for all neighboring databases X, X' and all $\alpha \in (1, \infty)$,*

$$D_\alpha(M(X) \| M(X')) \leq \rho\alpha,$$

87 *where D_α is the Renyi divergence [Van Erven and Harremoës, 2014].*

88 **Lemma 2.7** (Adaptive composition and Post Processing of zCDP [Bun and Steinke, 2016]). *Let*
 89 *$M : \mathcal{X}^n \rightarrow \mathcal{Y}$ and $M' : \mathcal{X}^n \times \mathcal{Y} \rightarrow \mathcal{Z}$. Suppose M satisfies ρ -zCDP and M' satisfies ρ' -zCDP*
 90 *(as a function of its first argument). Define $M'' : \mathcal{X}^n \rightarrow \mathcal{Z}$ by $M''(x) = M'(x, M(x))$. Then M''*
 91 *satisfies $(\rho + \rho')$ -zCDP.*

92 **Lemma 2.8** (Privacy Guarantee of Gaussian mechanism [Dwork et al., 2014, Bun and Steinke,
 93 2016]). *Let $f : \mathbb{N}^{\mathcal{X}} \mapsto \mathbb{R}^d$ be an arbitrary d -dimensional function with ℓ_2 sensitivity $\Delta_2 =$*
 94 *$\sup_{\text{neighboring } X, X'} \|f(X) - f(X')\|_2$. Then for any $\rho > 0$, Gaussian Mechanism with parameter*
 95 *$\sigma^2 = \frac{\Delta_2^2}{2\rho}$ satisfies ρ -zCDP.*

96 **Lemma 2.9** (Converting zCDP to DP [Bun and Steinke, 2016]). *If M satisfies ρ -zCDP then M*
 97 *satisfies $(\rho + 2\sqrt{\rho \log(1/\delta)}, \delta)$ -DP.*

98 As we use exclusively Gaussian mechanisms and their composition in our proposed algorithms,
 99 our method actually satisfies the stronger (ϵ, δ) -DP guarantees than what is implied by zCDP via
 100 techniques from [Balle and Wang, 2018, Dong et al., 2019], which reduces the ϵ parameter by a
 101 sizable fraction in typical parameter regimes. We stick to zCDP for clarity and generality, because all
 102 our results would apply without changes if we modify the noise into other mechanisms satisfying
 103 zCDP, e.g., the Discrete Gaussian Mechanism [Canonne et al., 2020].

104 2.2 Revisiting Linear Sketches

105 Charikar et al. [2002] proposed the **Count-Median** sketch, a randomized algorithm that summarizes
 106 a database and solves the frequency estimation problem. The Count-Median sketch uses a $d \times w$ array
 107 of counters, i.e., $C[d, w]$, where all the counters are initialized to **zero**, and has two sets of independent
 108 hash functions h and g . For each row r , the hash function h_r maps input items uniformly onto
 109 $\{1, \dots, w\}$ and the hash function g_r maps input items uniformly onto $\{-1, +1\}$. For item x with
 110 value $v \in \{-1, +1\}$, Count-Median sketch updates d counters, one per each row, based on the hash
 111 values such that for a particular row r , $g_r(x)$ will be added or subtracted to the counter at the $h_r(x)^{th}$
 112 index depending on whether x is being inserted or deleted respectively, as shown in Algorithm 1.
 113 Hence, the update time is $O(d)$. To estimate the frequency of item x , Count-Median sketch will output
 114 the median $_{1 \leq r \leq d} g_r(x) \cdot C[r, h_r(x)]$, as shown in Algorithm 2. By updating each row's counter
 115 based on the hashed value of either 1 or -1 and reporting the median for query, Count-Median
 116 sketch provides an unbiased estimate. To reduce the failure probability of bad estimations, d is set to
 117 $O(\log(1/\beta))$ and it uses $O(\frac{1}{\gamma} \log(\frac{1}{\beta}))$ space to solve the frequency estimation problem.

Algorithm 1 Linear Sketch Update(x, v)

1: **Input:** Item x with value $v \in \{-1, +1\}$, counter arrays C , and two sets of hash functions $\{h_1, \dots, h_{C.rows}\}$ and $\{g_1, \dots, g_{C.rows}\}$.
2: **for** $r \leftarrow 1, 2, \dots, C.rows$ **do**
3: $C[r, h_r(x)] \leftarrow C[r, h_r(x)] + v \cdot g_r(x)$
4: **end for**
5: **Output:** C .

118 Cormode and Muthukrishnan [2005] proposed the **Count-Min** sketch that shares the same initialization, update, and data structure as Count-Median sketch. Count-Min sketch also uses $O(\frac{1}{\gamma} \log(\frac{1}{\beta}))$
119 space to solve the frequency estimation problem. A major difference is that Count-Min sketch
120 makes all hash functions in set g return positive 1. As a result, to estimate the frequency of item x ,
121 Count-Min sketch returns $\min_{1 \leq r \leq d} C[r, h_r(x)]$ instead of the median, as shown in Algorithm 2. In
122 addition, it has the nice property of never underestimating item's frequency. Since linear sketches can
123 approximate an item's frequency accurately, they also solves the top K approximation problem by
124 returning the K items associated with the highest estimated frequency.
125

Algorithm 2 Linear Sketch Query(x)

1: **Input:** Item x , counter arrays C , and two sets of hash functions $\{h_1, \dots, h_{C.rows}\}$ and $\{g_1, \dots, g_{C.rows}\}$.
2: $arr \leftarrow []$
3: **for** $r \leftarrow 1, 2, \dots, C.rows$ **do**
4: $arr.append(g_r(x) \cdot C[r, h_r(x)])$
5: **end for**
6: **Output:** $\min(arr)$ **for Count-Min** or $\text{median}(arr)$ **for Count-Median**.

126 Moreover, Gilbert et al. [2002] made the connection between frequency and quantiles to propose
127 the first universe based *RSS* quantile sketch in the turnstile model. The observation is that the
128 quantile range, from 0 to the item, can be decomposed into at most $\log U$ dyadic intervals [Cormode
129 et al., 2019] and the sum over the estimated frequencies for these intervals gives the estimated rank.
130 Cormode and Muthukrishnan [2005] proposed the Dyadic Count-Min sketch (DCM) which uses
131 Count-Min sketches for estimating the frequencies of each dyadic interval with space complexity
132 $O(\frac{1}{\gamma} \log^2 U \log \frac{\log U}{\gamma})$ and update time $O(\log U \log \frac{\log U}{\gamma})$. Later, Wang et al. [2013] leveraged the
133 unbiased property of Count-Median sketches and proposed the Dyadic Count-Median sketch (DCS)
134 which replaces the Count-Min sketch with the Count-Median Sketch [Charikar et al., 2002] to further
135 improve the space complexity to $O(\frac{1}{\gamma} \log^{1.5} U \log^{1.5}(\frac{\log U}{\gamma}))$ while maintaining the same update time.
136 DCS and DCM share the same update and query algorithms as shown in Appendix B. DCM uses
137 $O(\frac{1}{\gamma} \log U \log \frac{\log U}{\gamma})$ space for each Count-Min sketch and DCS uses $O(\frac{1}{\gamma} \log^{0.5} U \log^{1.5}(\frac{\log U}{\gamma}))$
138 space for each Count-Median sketch, where both of them use $O(\log \frac{\log U}{\gamma})$ rows. For more specific
139 details, [Cormode and Yi, 2020] provide a comprehensive analysis of linear and quantile sketches.

140 3 Private Linear Sketches

141 In this section, we present new algorithms for differentially private linear sketches. We highlight
142 that our Private Count Min/Median only require a different initialization while they share the same
143 update (Algorithm 1) and query (Algorithm 2) with the original Count Min/Median. Therefore, the
144 implementation of our algorithms is efficient. Below we show our private initialization.

Algorithm 3 DP Linear Sketch Initialization with Gaussian Noise

```
1: Input: Desired accuracy parameter  $\gamma$ , failure probability  $\beta$ , and budget for zCDP  $\rho$ .
2: Initialize Counter Arrays
3:  $\sigma \leftarrow \sqrt{\log(2/\beta)/\rho}$ 
4:  $E \leftarrow \sqrt{\frac{2 \log \frac{2}{\beta}}{\rho}} \cdot \sqrt{\log \frac{\frac{4}{\gamma} \log(\frac{2}{\beta})}{\beta}}$ 
5: for  $r \leftarrow 1, 2, \dots, \log(2/\beta)$  do
6:   for  $c \leftarrow 1, 2, \dots, 1/\gamma$  do
7:      $C[r, c] \leftarrow \mathcal{N}(0, \sigma^2)$  if Private Count-Median
8:      $C[r, c] \leftarrow E + \mathcal{N}(0, \sigma^2)$  if Private Count-Min
9:   end for
10: end for
11: Output:  $C$ .
```

145 In Algorithm 3, the set of arrays we use is C which consists of $\log(2/\beta)$ arrays with length $1/\gamma$,
146 which has the same space complexity as original Count Min/Median. Recall that two neighboring
147 databases X and X' differ by at most one item. Therefore, after updating all the items respectively,
148 for each corresponding array in $C(X)$ and $C(X')$, they differ by at most two elements and the
149 difference is at most 1. Then the ℓ_2 -sensitivity of the set of arrays C is bounded by

$$\Delta_2 = \sqrt{2 \log(2/\beta)}. \quad (1)$$

150 By applying the Gaussian Mechanism (Definition 2.5), we can add *independent* Gaussian noises
151 $\mathcal{N}(0, \sigma^2)$ to each counter in C , where $\sigma = \sqrt{\frac{\log(2/\beta)}{\rho}}$. Due to the privacy guarantee of Gaussian
152 Mechanism (Lemma 2.8), it satisfies $\frac{\Delta_2^2}{2\sigma^2} = \rho$ -zCDP.

153 Define $E(\beta, \gamma, \rho) = \sqrt{\frac{2 \log \frac{2}{\beta}}{\rho}} \cdot \sqrt{\log \frac{\frac{4}{\gamma} \log(\frac{2}{\beta})}{\beta}}$, for simplicity, we will use E in Algorithm 3 and
154 the proof in Appendix A. The private version of Count Min can be derived by adding *independent*
155 Gaussian noises $\mathcal{N}(E, \sigma^2)$ to each counter of C , while the private version of Count Median can be
156 derived by adding *independent* Gaussian noises $\mathcal{N}(0, \sigma^2)$ to each counter of C . The private version
157 of Count Min/Median is derived by combining Algorithm 3, Algorithm 1, and Algorithm 2.

158 **Some explanations.** Note that here the difference in the noises we add to Private Count Min/Median
159 is because we want to retain the nice properties of Count Min (never underestimating) and Count
160 Median (unbiased estimate).

161 3.1 Main results about Private Count Min/Median

162 We present the privacy guarantee and utility analysis of our Private Count Min/Median below. Recall
163 that for each item x , we perform update as in Algorithm 1 and query as in Algorithm 2. In addition,
164 $\hat{f}(x)$ is the output estimated frequency and $f(x)$ is the actual frequency. We begin with the properties
165 of Private Count Min. Note that all the proofs are deferred to Appendix A.

166 **Theorem 3.1.** *Private Count Min satisfies ρ -zCDP regardless of the number of queries. Furthermore,*
167 *for each item x , with probability $1 - \beta$, the output $\hat{f}(x)$ satisfies that*

$$0 \leq \hat{f}(x) - f(x) \leq \gamma \cdot N + 2E = \gamma \cdot N + 2\sqrt{\frac{2 \log \frac{2}{\beta}}{\rho}} \cdot \sqrt{\log \frac{\frac{4}{\gamma} \log(\frac{2}{\beta})}{\beta}}. \quad (2)$$

168 **Comparison to Count Min.** Comparing our Theorem 3.1 with the conclusion in [Cormode and
169 Muthukrishnan, 2005], our Private Count Min preserves the nice property that the output will not
170 underestimate the frequency with high probability. Furthermore, within the most popular regime
171 where the privacy budget ρ is a constant, the additional error bound due to differential privacy is
172 independent of the size of database N , therefore it will become negligible as N goes large.

173 **Justification of our Gaussian noise.** Note that with high probability, all the noises we add ($E + \sigma_{i,j}$,
174 $\sigma_{i,j} \sim \mathcal{N}(0, \sigma^2)$) will be non-negative. Therefore, the noise we add and the original error induced by

Count Min will directly sum up and lead to larger error in evaluation. However, we claim that the additional E ensures that with high probability, for all item x , the output will not underestimate the actual frequency. This nice property enables the good performance of our Private Count Min when used in approximate top k task, as shown in Section 5.

Next, Theorem 3.2 shows the properties of Private Count Median.

Theorem 3.2. *Private Count Median satisfies ρ -zCDP regardless of the number of queries. Furthermore, for each item x , the output $\hat{f}(x)$ satisfies that $\mathbb{E}\hat{f}(x) = f(x)$ and with probability $1 - \beta$,*

$$|\hat{f}(x) - f(x)| \leq \gamma \cdot N + E = \gamma \cdot N + \sqrt{\frac{2 \log \frac{2}{\beta}}{\rho}} \cdot \sqrt{\log \frac{\frac{4}{\gamma} \log(\frac{2}{\beta})}{\beta}}. \quad (3)$$

Comparison to Count Median. Comparing our Theorem 3.2 with the conclusion in [Charikar et al., 2002], our Private Count Median preserves the nice property that the output will be an unbiased estimate of the frequency. This property enables our use of Private Count Median in quantile estimation below (Section 4). Furthermore, within the most popular regime where the privacy budget ρ is a constant, the additional error bound due to differential privacy is independent of the size of database N , thus it will become negligible as N goes large.

Better performance in evaluation. Since the noise we add is independent to the error of Count Median, these two errors may cancel each other out. In addition, the scale of the whole error is smaller than a direct summation of the scale of these two errors, which will lead to a much better performance in evaluation (Section 5) compared to the worst case bound.

4 Private Quantile Sketches

In this section, we apply our Private Count Median to state of the art quantile sketches in the turnstile model. Our private Dyadic Count-Median sketch can estimate all the quantiles accurately at the same time while ensuring differential privacy.

4.1 Revisiting DCS

In [Wang et al., 2013], it is shown that DCS can return all γ -approximate quantiles with constant probability using space $O\left(\frac{1}{\gamma} \log^{1.5} U \log^{1.5}\left(\frac{\log U}{\gamma}\right)\right)$. More specifically, the sketch structure here consists of $\log U$ Count Medians, each Count Median uses a counter arrays C , which is $d \times w$ counters. The choice of d, w follows $d = \Theta\left(\log\left(\frac{\log U}{\gamma}\right)\right)$ and $w = O\left(\sqrt{\log U \log\left(\frac{\log U}{\gamma}\right)}/\gamma\right)$.

4.2 Private DCS

In this work, we aim to estimate the quantiles accurately while preserving privacy. We do this by replacing Count Median with Private Count Median, which bases on the same structure as Count Median discussed above. Given the privacy budget ρ , the privacy budget of each Private Count Median is thus $\rho_0 = \frac{\rho}{\log U}$, due to composition of zCDP (Lemma 2.7). The ℓ_2 -sensitivity of each Private Count Median is

$$\Delta_2 = O(\sqrt{2d}) = O\left(\sqrt{\log\left(\frac{\log U}{\gamma}\right)}\right).$$

To keep the whole algorithm ρ -zCDP, it suffices to keep each Count Median ρ_0 -zCDP (Lemma 2.7). Therefore, Gaussian Mechanism (Definition 2.5) with $\sigma^2 = O\left(\log U \log\left(\frac{\log U}{\gamma}\right)/\rho\right)$ ensures ρ -zCDP (Lemma 2.8). Similar to Lemma A.1, define $E(\gamma, U) = O\left(\sqrt{\frac{\log U \log\left(\frac{\log U}{\gamma}\right)}{\rho}} \cdot \sqrt{\log \frac{\log U \log\left(\frac{\log U}{\gamma}\right)}{\gamma}}\right)$, we can prove that with constant probability, all the

207 Gaussian noises we add to all $\log U$ Count Medians are bounded by E (for simplicity, we use E to
208 represent $E(\gamma, U)$).

Conditioned on the high probability event above, we prove that for a fixed quantile, the estimated quantile will be accurate with high probability. As has been proven in Theorem 3.2, the output estimated frequency is unbiased for any item. Therefore, similar to [Wang et al., 2013], for any item x (corresponding to a fixed Count Median), we have the output $\hat{f}(x)$ of that Count Median satisfies

$$\mathbb{P}\left[\left|\hat{f}(x) - f(x)\right| > \frac{1}{w} \cdot N + E\right] < \exp(-O(d)) = O\left(\frac{\gamma}{\log U}\right).$$

209 By a union bound, with probability $1 - \log U \times O\left(\frac{\gamma}{\log U}\right) = 1 - O(\gamma)$, for any item corresponding
210 to this fixed quantile, the error of Count Median is bounded by $\frac{1}{w} \cdot N + E$. Conditioned upon this
211 event, by Hoeffding's inequality, with probability $1 - O\left(\frac{\gamma}{\log U}\right)$, the sum of $\log U$ such independent
212 errors is bounded by

$$\sqrt{\log U \log\left(\frac{\log U}{\gamma}\right)} \cdot \left(\frac{N}{w} + E\right) = \gamma \cdot N + E', \quad (4)$$

213 where $E' = O\left(\frac{\log U \log\left(\frac{\log U}{\gamma}\right)}{\sqrt{\rho}} \cdot \sqrt{\log\left(\frac{\log U \log\left(\frac{\log U}{\gamma}\right)}{\gamma}\right)}\right)$. To sum up, for a fixed quantile, with
214 probability $1 - O(\gamma)$, the estimating error is bounded by $\gamma \cdot N + E'$.

215 Finally, apply another union bound on the $\frac{1}{\gamma}$ different quantiles, with constant probability, all the
216 quantiles are estimated accurately (within the error bound (4)). Note that similar to [Wang et al.,
217 2013], the failure probability here is a constant. For any failure probability β , we can further increase
218 d by a factor of $\log \frac{1}{\beta}$ to reduce this failure probability to β .

219 **Take-away of Private DCS.** First, our Private DCS has a same space complexity as the original
220 DCS. In addition, according to (4), the additional error bound is proportional to $\frac{\log U \log \frac{1}{\gamma}}{\sqrt{\rho}}$ (ignoring
221 $\log \log$ terms), and independent to the size of database N . In the most popular regime where the
222 privacy budget ρ is a constant, the additional error bound only appears as lower order terms, which
223 will become negligible as N goes large.

224 5 Evaluation

225 5.1 Data Sets

226 The experimental evaluation is conducted using both synthetic and real world data sets. We con-
227 sider the synthetic Zipf dataset Zipf [2016] and the real CAIDA Anonymized Internet Trace 2015
228 dataset pas. For each independent run in the experiments, we use an input database size $N = 10^5$.

- 229 • **Zipf Dataset:** Items are drawn from a bounded universe of size 2^{16} and items' frequencies
230 follow the Zipf Law Zipf [2016].
- 231 • **2015 CAIDA Dataset:** The CAIDA dataset is collected from the Chicago high-speed
232 monitor. Items are the source IP addresses with bounded universe size of 2^{32} .

233 5.2 Metrics

234 In all experiments, we average the various metrics over 5 independent runs to minimize the measure-
235 ment variance as all data sketches have randomness depending on the hash function. The metrics
236 used in the experiments are:

- 237 • **Average Relative Error:** Let the set Ψ denotes all unique items in the database. The average
238 relative error (ARE) is computed based on Ψ in which $\frac{1}{\Psi} \sum_{e \in \Psi} \frac{|f(e) - \hat{f}(e)|}{f(e)}$.

- **F1 Score:** F1 score is the harmonic mean of the precision and recall ($2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$).
- **Average Rank Error:** For each evenly spaced quantile and its associated item, we average the distance between the true rank and estimated rank.

We use ARE to evaluate sketch performance on frequency estimation and F1 score to evaluate the sketch’s performance in identifying the top 10 items. For quantile approximation, we consider the m evenly spaced quantiles and items. For instance, if $m = 1$, we consider the rank error for the median item; if $m = 2$, we consider then average rank error for the 33rd and 67th percentile items. Lower ARE and average rank error, and higher F1 score indicate better approximation.

5.3 Private Linear Sketches Experiments

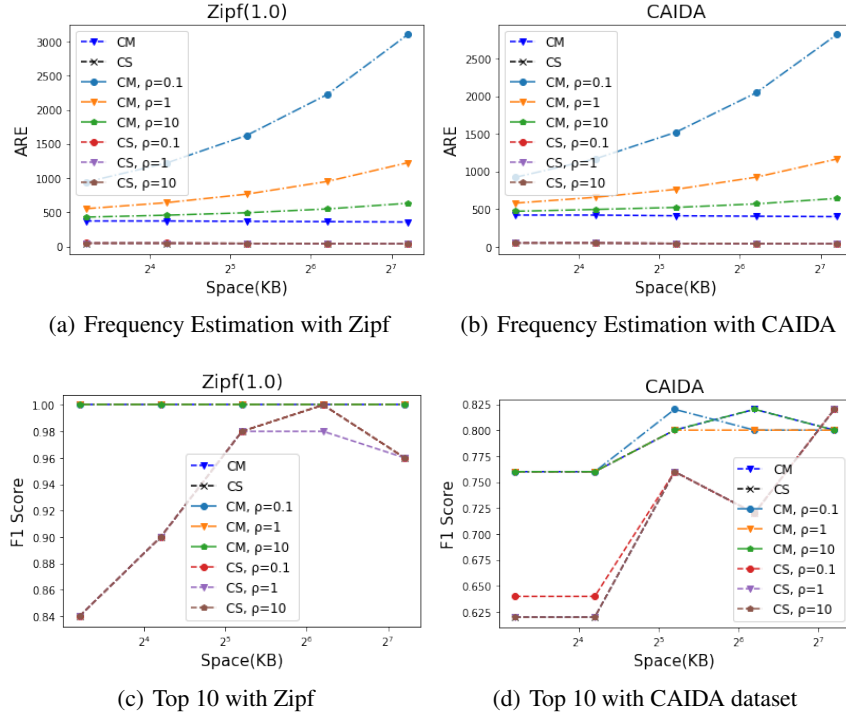


Figure 1: Comparison of non-private linear sketches and DP linear sketches with various privacy budget under synthetic and real world datasets. The experiments assume $\beta = 1\%$ and $N = 10^5$.

To evaluate the utility of DP linear sketches, we compare the average relative error (ARE) and F1 score for frequency estimation and identify the top 10 items, respectively. As shown in Figure 1, the x-axis represents the space budget for each sketch (from 9.2 KB to 147.3 KB), and the y-axis denotes ARE or F1 score. The DP linear sketches use $\rho \in \{0.1, 1, 10\}$ in which lower ρ value indicates more noise need to be added, and all sketches assume $\beta = 1\%$.

For frequency estimation, the performance of private Count-Median sketches with various privacy budgets is basically equivalent to the performance of the non-private Count-Median sketch. Under different space and privacy budgets, they have minimal difference in ARE for both Zipf and CAIDA datasets, meaning that, while providing strong privacy guarantee, the estimated frequencies are still very accurate. The accurate estimation of private Count-Median is primarily due to the unbiased nature of Count-Median in which, by adding Gaussian noise, the private Count-Median still provides unbiased estimation for an item’s frequency as proved in Theorem 3.2. As shown in both Figure 1(a) and Figure 1(b), the performance of private Count-Min degrades when the space budget increases or the privacy budget decreases. This behavior is expected as the upper bound on the frequency

error in Theorem 3.1 has a dependency on both γ and ρ . In order to preserve the property of not underestimating an item's frequency, the private Count-Min sketch needs to add larger noise to each counter when the number of counters increases. As a result, the estimated frequencies for low-frequency items become inflated and this in turn decreases the overall accuracy.

For approximate top 10 items, private Count-Median has similar performance to Count-Median. Since non-private and private Count-Median are unbiased, they may underestimate the frequency of true top K items and decrease the recall. On the other hand, the property of no underestimation is desirable for approximate top K items. In particular, non-private and private Count-Min sketch score high F1 scores for all datasets. While providing privacy guarantees, private Count-Min achieve 1.0 F1 scores for all space and all privacy budgets in Zipf dataset, as shown in Figure 1(c).

5.4 Private Quantile Sketch Experiments

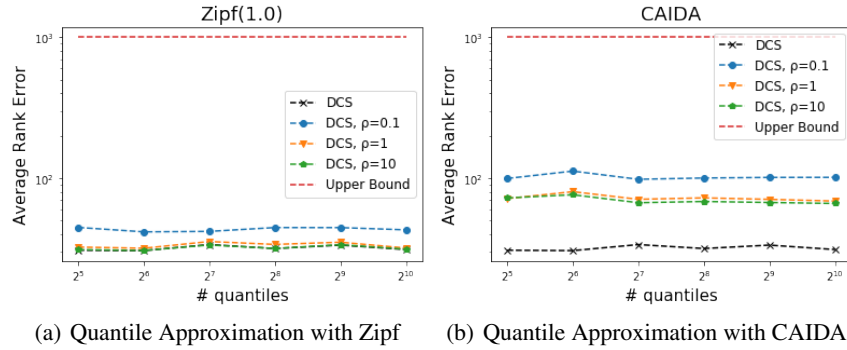


Figure 2: Compare DCS and DP DCS with various privacy budget under synthetic and real world datasets. The experiments assume $\gamma = 1\%$, $N = 10^5$, and the desired error upper bound is 10^3 (γN).

To evaluate the utility of DP DCS, we compare the average rank error. As shown in Figure 2, the x-axis represents the number of evenly spaced quantiles, and the y-axis denotes the average rank error. The DP DCS use privacy budget $\rho \in \{0.1, 1, 10\}$ and all sketches assume $\gamma = 1\%$.

For the quantile approximation, we observe that the increase in the number of evenly spaced quantiles does not impact the average rank error, as shown in both Figure 2(a) and Figure 2(b). Since the CAIDA dataset universe size (2^{32}) is larger than Zipf dataset universe size (2^{16}), the average rank error in the CAIDA dataset is larger than the average rank error in the Zipf dataset. As shown in Equation (4), the error bound has a term depending on the universe size in which a large universe size leads to more error. When the privacy budget decreases, the average rank error increases as more noise needs to be added. Comparing DP DCS with strong privacy ($\rho = 0.1$) to DCS, the increase in rank error is relatively small compared to the database size of 10^5 . In addition, the desired rank error upper bound is $\gamma \cdot N = 10^3$ and all the rank errors are one order of magnitude lower.

6 Conclusion

In this work, we demonstrate that linear sketches can be made differentially private and provide useful information while maintaining their original properties by adding a small amount of Gaussian noise at initialization. In addition, leveraging the private Count-Median sketch, we propose the DP DCS for quantile approximation in the turnstile model, and the DP DCS achieves low rank error even for a large data universe. Moreover, for all proposed algorithms, when the privacy budget is constant, the additional error due to privacy is independent of the database size and the error will become negligible when the database grows larger. As a result, we believe our proposed algorithms are efficient and practical for real-world systems and enable these systems to perform data analysis and machine learning tasks privately.

References

- Anonymized internet traces 2015. https://catalog.caida.org/details/dataset/passive_2015_pcap. Accessed: 2022-5-10.
- Daniel Alabi, Omri Ben-Eliezer, and Anamay Chaturvedi. Bounded space differentially private quantiles. *arXiv preprint arXiv:2201.03380*, 2022.
- Raman Arora, Jalaj Upadhyay, et al. Differentially private robust low-rank approximation. *Advances in neural information processing systems*, 31, 2018.
- Peter Bailis, Edward Gan, Samuel Madden, Deepak Narayanan, Kexin Rong, and Sahaana Suri. Macrobase: Prioritizing attention in fast data. In *Proceedings of the 2017 ACM International Conference on Management of Data*, pages 541–556, 2017.
- Borja Balle and Yu-Xiang Wang. Improving the gaussian mechanism for differential privacy: Analytical calibration and optimal denoising. In *International Conference on Machine Learning*, pages 394–403. PMLR, 2018.
- Jeremiah Blocki, Avrim Blum, Anupam Datta, and Or Sheffet. The johnson-lindenstrauss transform itself preserves differential privacy. In *2012 IEEE 53rd Annual Symposium on Foundations of Computer Science*, pages 410–419. IEEE, 2012.
- Mark Bun and Thomas Steinke. Concentrated differential privacy: Simplifications, extensions, and lower bounds. In *Theory of Cryptography Conference*, pages 635–658. Springer, 2016.
- Clément L Canonne, Gautam Kamath, and Thomas Steinke. The discrete gaussian for differential privacy. *Advances in Neural Information Processing Systems*, 33:15676–15688, 2020.
- Moses Charikar, Kevin Chen, and Martin Farach-Colton. Finding frequent items in data streams. In *International Colloquium on Automata, Languages, and Programming*, pages 693–703. Springer, 2002.
- Seung Geol Choi, Dana Dachman-Soled, Mukul Kulkarni, and Arkady Yerukhimovich. Differentially-private multi-party sketching for large-scale statistics. *Cryptology ePrint Archive*, 2020.
- Graham Cormode. Current trends in data summaries. *ACM SIGMOD Record*, 50(4):6–15, 2022.
- Graham Cormode and Marios Hadjieleftheriou. Finding frequent items in data streams. *Proceedings of the VLDB Endowment*, 1(2):1530–1541, 2008.
- Graham Cormode and Shan Muthukrishnan. An improved data stream summary: the count-min sketch and its applications. *Journal of Algorithms*, 55(1):58–75, 2005.
- Graham Cormode and Ke Yi. *Small summaries for big data*. Cambridge University Press, 2020.
- Graham Cormode, Tejas Kulkarni, and Divesh Srivastava. Answering range queries under local differential privacy. *Proceedings of the VLDB Endowment*, 12(10):1126–1138, 2019.
- Sudipto Das, Shyam Antony, Divyakant Agrawal, and Amr El Abbadi. Thread cooperation in multicore architectures for frequency counting over multiple data streams. *Proceedings of the VLDB Endowment*, 2(1):217–228, 2009.
- Charlie Dickens, Justin Thaler, and Daniel Ting. (nearly) all cardinality estimators are differentially private. *arXiv preprint arXiv:2203.15400*, 2022.
- Jinshuo Dong, Aaron Roth, and Weijie J Su. Gaussian differential privacy. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 2019.
- Cynthia Dwork and Guy N Rothblum. Concentrated differential privacy. *arXiv preprint arXiv:1603.01887*, 2016.

337 Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in
338 private data analysis. In *Theory of cryptography conference*, pages 265–284. Springer, 2006.

339 Cynthia Dwork, Aaron Roth, et al. The algorithmic foundations of differential privacy. *Found. Trends*
340 *Theor. Comput. Sci.*, 9(3-4):211–407, 2014.

341 Anna C Gilbert, Yannis Kotidis, S Muthukrishnan, and Martin J Strauss. How to summarize the
342 universe: Dynamic maintenance of quantiles. In *VLDB’02: Proceedings of the 28th International*
343 *Conference on Very Large Databases*, pages 454–465. Elsevier, 2002.

344 Jennifer Gillenwater, Matthew Joseph, and Alex Kulesza. Differentially private quantiles. In
345 *International Conference on Machine Learning*, pages 3713–3722. PMLR, 2021.

346 Michael Greenwald and Sanjeev Khanna. Space-efficient online computation of quantile summaries.
347 *ACM SIGMOD Record*, 30(2):58–66, 2001.

348 Arpit Gupta, Rüdiger Birkner, Marco Canini, Nick Feamster, Chris Mac-Stoker, and Walter Willinger.
349 Network monitoring as a streaming analytics problem. In *Proceedings of the 15th ACM workshop*
350 *on hot topics in networks*, pages 106–112, 2016.

351 Nikita Ivkin, Zhuolong Yu, Vladimir Braverman, and Xin Jin. Qpipe: Quantiles sketch fully in
352 the data plane. In *Proceedings of the 15th International Conference on Emerging Networking*
353 *Experiments And Technologies*, pages 285–291, 2019.

354 Rajesh Jayaram and David P Woodruff. Data streams with bounded deletions. In *Proceedings of*
355 *the 37th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*, pages
356 341–354, 2018.

357 Zohar Karnin, Kevin Lang, and Edo Liberty. Optimal quantile approximation in streams. In *2016*
358 *IEEE 57th annual symposium on foundations of computer science (focs)*, pages 71–78. IEEE, 2016.

359 Tian Li, Zaoxing Liu, Vyas Sekar, and Virginia Smith. Privacy for free: Communication-efficient
360 learning with differential privacy using sketches. *arXiv preprint arXiv:1911.00972*, 2019.

361 Ahmed Metwally, Divyakant Agrawal, and Amr El Abbadi. Efficient computation of frequent and
362 top-k elements in data streams. In *International conference on database theory*, pages 398–412.
363 Springer, 2005.

364 Darakhshan Mir, Shan Muthukrishnan, Aleksandar Nikolov, and Rebecca N Wright. Pan-private
365 algorithms via statistics on sketches. In *Proceedings of the thirtieth ACM SIGMOD-SIGACT-*
366 *SIGART symposium on Principles of database systems*, pages 37–48, 2011.

367 Jayadev Misra and David Gries. Finding repeated elements. *Science of computer programming*, 2(2):
368 143–152, 1982.

369 J Ian Munro and Mike S Paterson. Selection and sorting with limited storage. *Theoretical computer*
370 *science*, 12(3):315–323, 1980.

371 Rasmus Pagh and Nina Mesing Stausholm. Efficient differentially private f_0 linear sketching. *arXiv*
372 *preprint arXiv:2001.11932*, 2020.

373 Daniel Rothchild, Ashwinee Panda, Enayat Ullah, Nikita Ivkin, Ion Stoica, Vladimir Braverman,
374 Joseph Gonzalez, and Raman Arora. Fetchsgd: Communication-efficient federated learning with
375 sketching. In *International Conference on Machine Learning*, pages 8253–8265. PMLR, 2020.

376 Nisheeth Shrivastava, Chiranjeev Buragohain, Divyakant Agrawal, and Subhash Suri. Medians and
377 beyond: new aggregation techniques for sensor networks. In *Proceedings of the 2nd international*
378 *conference on Embedded networked sensor systems*, pages 239–249, 2004.

- Adam Smith, Shuang Song, and Abhradeep Guha Thakurta. The flajolet-martin sketch itself preserves differential privacy: Private counting with minimal space. *Advances in Neural Information Processing Systems*, 33:19561–19572, 2020.
- Christos Tzamos, Emmanouil-Vasileios Vlatakis-Gkaragkounis, and Ilias Zadik. Optimal private median estimation under minimal distributional assumptions. *Advances in Neural Information Processing Systems*, 33:3301–3311, 2020.
- Jalaj Upadhyay. Differentially private linear algebra in the streaming model. *arXiv preprint arXiv:1409.5414*, 2014.
- Salil Vadhan. The complexity of differential privacy. In *Tutorials on the Foundations of Cryptography*, pages 347–450. Springer, 2017.
- Tim Van Erven and Peter Harremoës. Rényi divergence and kullback-leibler divergence. *IEEE Transactions on Information Theory*, 60(7):3797–3820, 2014.
- Lu Wang, Ge Luo, Ke Yi, and Graham Cormode. Quantiles over data streams: an experimental study. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, pages 737–748, 2013.
- Victor Zakhary, Lawrence Lim, Divyakant Agrawal, and Amr El Abbadi. Cot: Decentralized elastic caches for cloud environments. *arXiv preprint arXiv:2006.08067*, 2020.
- Fuheng Zhao, Divyakant Agrawal, Amr El Abbadi, and Ahmed Metwally. Spacesaving \pm : An optimal algorithm for frequency estimation and frequent items in the bounded deletion model. *arXiv preprint arXiv:2112.03462*, 2021a.
- Fuheng Zhao, Sujaya Maiyya, Ryan Wiener, Divyakant Agrawal, and Amr El Abbadi. Kll \pm : Approximate quantile sketches over dynamic datasets. *Proceedings of the VLDB Endowment*, 14(7):1215–1227, 2021b.
- George Kingsley Zipf. *Human behavior and the principle of least effort: An introduction to human ecology*. Ravenio Books, 2016.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? **[Yes]** Claims in abstract and introduction reflect our contributions.
 - (b) Did you describe the limitations of your work? **[Yes]** In Theorem 3.1, the private Count-Min frequency estimation’s additional error has a dependency of $\sqrt{\log \frac{1}{\gamma}}$ in order to keep the property of not underestimating item’s frequency. Assume database size is fixed, decreasing gamma will increase the average error. Note the error is independent from the database size, and when database size grows, the error will become negligible.
 - (c) Did you discuss any potential negative societal impacts of your work? **[N/A]** Apply our proposed algorithms to current systems will protect user privacy.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**
 - (b) Did you include complete proofs of all theoretical results? **[Yes]** We put some proves in Appendix A due to space limits

- 422 3. If you ran experiments...
- 423 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
- 424 mental results (either in the supplemental material or as a URL)? [Yes] We will provide
- 425 an url to our Github Repo once the paper is accepted.
- 426 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
- 427 were chosen)? [Yes] We provide the parameters used for the experiments.
- 428 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
- 429 ments multiple times)? [N/A] We didn't use random seed
- 430 (d) Did you include the total amount of compute and the type of resources used (e.g.,
- 431 type of GPUs, internal cluster, or cloud provider)? [N/A] we didn't use any external
- 432 resources beside a macbook pro.
- 433 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 434 (a) If your work uses existing assets, did you cite the creators? [N/A] We did not use
- 435 existing assets.
- 436 (b) Did you mention the license of the assets? [N/A]
- 437 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 438
- 439 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 440 using/curating? [N/A]
- 441 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 442 information or offensive content? [N/A]
- 443 5. If you used crowdsourcing or conducted research with human subjects...
- 444 (a) Did you include the full text of instructions given to participants and screenshots, if
- 445 applicable? [N/A] We did not use crowdsourcing or conducted research with human
- 446 subjects.
- 447 (b) Did you describe any potential participant risks, with links to Institutional Review
- 448 Board (IRB) approvals, if applicable? [N/A]
- 449 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 450 spent on participant compensation? [N/A]

A Missing proofs

In this section, we present the missing proofs. Recall that for each item x , we perform query as in Algorithm 2. In the proof below, we use arr' to denote the arr in Algorithm 2 (with private initialization) while arr denotes the arr in Algorithm 2 under non-private initialization (all zero initialization) and the same set of hash functions. In addition, $\hat{f}(x)$ is the output estimated frequency and $f(x)$ is the actual frequency. We first present the following Lemma A.1, which gives a high probability bound for the Gaussian noises we add.

Lemma A.1 (Utility analysis). *If there are $\frac{1}{\gamma} \log(\frac{2}{\beta})$ independent Gaussian noises sampled from $\mathcal{N}(0, \sigma^2)$ (where $\sigma = \sqrt{\frac{\log(2/\beta)}{\rho}}$), denoted as σ_{ij} where $i \in [\log(2/\beta)]$, $j \in [1/\gamma]$, then with probability $1 - \frac{\beta}{2}$, for any $i \in [\log(2/\beta)]$, $j \in [1/\gamma]$,*

$$|\sigma_{ij}| \leq \sqrt{2}\sigma \cdot \sqrt{\log \frac{\frac{4}{\gamma} \log(\frac{2}{\beta})}{\beta}} = \sqrt{\frac{2 \log \frac{2}{\beta}}{\rho}} \cdot \sqrt{\log \frac{\frac{4}{\gamma} \log(\frac{2}{\beta})}{\beta}}. \quad (5)$$

Proof of Lemma A.1. The lemma directly results from the concentration inequality of Gaussian distribution and a union bound. \square

Theorem A.2 (Restate Theorem 3.1). *Private Count Min satisfies ρ -zCDP regardless of the number of queries. Furthermore, for each item x , with probability $1 - \beta$, the output $\hat{f}(x)$ satisfies that*

$$0 \leq \hat{f}(x) - f(x) \leq \gamma \cdot N + 2E = \gamma \cdot N + 2\sqrt{\frac{2 \log \frac{2}{\beta}}{\rho}} \cdot \sqrt{\log \frac{\frac{4}{\gamma} \log(\frac{2}{\beta})}{\beta}}. \quad (6)$$

Proof of Theorem A.2. Differential privacy directly results from the DP guarantee of Gaussian Mechanism (Lemma 2.8) and post processing (Lemma 2.7).

By the property of Count Min [Cormode and Muthukrishnan, 2005], we have for any item x , with probability $1 - \frac{\beta}{2}$,

$$0 \leq \min(\text{arr}) - f(x) \leq \gamma \cdot N.$$

According to Lemma A.1, we have with probability $1 - \frac{\beta}{2}$, for all noises $E + \sigma_{ij}$, where $\sigma_{ij} \sim \mathcal{N}(0, \sigma^2)$, $i \in [\log(2/\beta)]$, $j \in [1/\gamma]$, it holds that

$$0 \leq E + \sigma_{ij} \leq 2E.$$

Conditioned on these two cases that will happen with probability $1 - \beta$, it holds that

$$\min(\text{arr}) \leq \min(\text{arr}') \leq \min(\text{arr} + 2E \cdot \mathbf{1}) \leq \min(\text{arr}) + 2E.$$

Therefore, we have $\hat{f}(x) = \min(\text{arr}') \geq \min(\text{arr}) \geq f(x)$ and

$$\hat{f}(x) = \min(\text{arr}') \leq \min(\text{arr}) + 2E \leq f(x) + \gamma \cdot N + 2E.$$

Then the proof is completed by plugging in the definition of E . \square

Theorem A.3 (Restate Theorem 3.2). *Private Count Median satisfies ρ -zCDP regardless of the number of queries. Furthermore, for each item x , the output $\hat{f}(x)$ satisfies that $\mathbb{E}\hat{f}(x) = f(x)$ and with probability $1 - \beta$,*

$$|\hat{f}(x) - f(x)| \leq \gamma \cdot N + E = \gamma \cdot N + \sqrt{\frac{2 \log \frac{2}{\beta}}{\rho}} \cdot \sqrt{\log \frac{\frac{4}{\gamma} \log(\frac{2}{\beta})}{\beta}}. \quad (7)$$

Proof of Theorem A.3. First of all, differential privacy directly results from the DP guarantee of Gaussian Mechanism (Lemma 2.8) and post processing (Lemma 2.7).

Next we claim that the conclusion that $\mathbb{E}\hat{f}(x) = f(x)$ arises from symmetry. If we replace $\{h_i\}_{i \in [\log(2/\beta)]}$, $\{g_i\}_{i \in [\log(2/\beta)]}$, $\{\sigma_{ij}\}_{i \in [\log(2/\beta)], j \in [1/\gamma]}$ with $\{h_i\}_{i \in [\log(2/\beta)]}$, $\{g'_i\}_{i \in [\log(2/\beta)]}$,

473 $\{-\sigma_{ij}\}_{i \in [\log(2/\beta)], j \in [1/\gamma]}$ where $g'_i(x) = g_i(x)$ and $g'_i(x') = -g_i(x')$, $\forall x' \neq x$, then the out-
474 puts under these two cases will be symmetric around $f(x)$ and the probability distribution function at
475 these two cases are identical. Therefore, we have

$$\mathbb{E}\hat{f}(x) = f(x). \quad (8)$$

Finally, by the property of Count Median [Charikar et al., 2002], we have for any item x , with probability $1 - \frac{\beta}{2}$,

$$|\text{median}(\text{arr}) - f(x)| \leq \gamma \cdot N.$$

According to Lemma A.1, we have with probability $1 - \frac{\beta}{2}$, for all noises $\sigma_{ij} \sim \mathcal{N}(0, \sigma^2)$, $i \in [\log(2/\beta)]$, $j \in [1/\gamma]$, it holds that

$$|\sigma_{ij}| \leq E.$$

Conditioned on these two cases that will happen with probability $1 - \beta$, we will prove that

$$|\hat{f}(x) - f(x)| \leq \gamma \cdot N + E.$$

Without loss of generality, we can assume that $\log \frac{2}{\beta} = 2k + 1$, where k is a positive integer (we can choose k to be the minimum integer such that $2k + 1 \geq \log \frac{2}{\beta}$).

Suppose that $\text{median}(\text{arr}') - f(x) > \gamma \cdot N + E$, then it holds that there are at least $k + 1$ elements in arr' that is larger than $f(x) + \gamma \cdot N + E$ due to the definition of median. Because $|\sigma_{ij}| \leq E$, for all i, j , there are at least $k + 1$ elements in arr that is larger than $f(x) + \gamma \cdot N$. Therefore, $\text{median}(\text{arr}) > f(x) + \gamma \cdot N$, which leads to contradiction. As a result, we have

$$\text{median}(\text{arr}') - f(x) \leq \gamma \cdot N + E.$$

Similarly, $\text{median}(\text{arr}') - f(x) \geq -\gamma \cdot N - E$. Combining these two results,

$$|\hat{f}(x) - f(x)| \leq \gamma \cdot N + E.$$

476 Then the proof is completed by plugging in the definition of E . □

477 B Missing Quantile Algorithms

Algorithm 4 DCS/DCM Update(x, v)

```

1: Input: Item  $x$  with value  $v \in \{-1, +1\}$ , and an array of linear sketches  $\{\text{LS}_0, \dots, \text{LS}_{\log U}\}$ .
2: for  $j \leftarrow 0, \dots, \log U$  do
3:    $\text{LS}_j.\text{update}(x, v)$ 
4:    $x \leftarrow \lfloor x/2 \rfloor$ 
5: end for
6: Output:  $\{\text{LS}_0, \dots, \text{LS}_{\log U}\}$ .
```

Algorithm 5 DCS/DCM Query(x)

```

1: Input: Item  $x$ , and an array of linear sketches  $\{\text{LS}_0, \dots, \text{LS}_{\log U}\}$ .
2:  $R \leftarrow 0$ 
3: for  $i \leftarrow 0, \dots, \log U$  do
4:   if  $x$  is odd then
5:      $R \leftarrow R + \text{LS}_i.\text{query}(x)$ 
6:   end if
7:    $x \leftarrow \lfloor x/2 \rfloor$ 
8: end for
9: Output:  $R$ .
```

478 Observe, that the relationship between frequency and rank is that one can sum up all items' frequency
479 in the range of 0 to the item itself to estimate the rank. However, this naive approach requires
480 summing all items' frequencies in the range, and the error quickly escalates. A better approach is to

481 break the range from 0 to item x into at most $\log U$ dyadic intervals and then sum all frequencies for
482 each dyadic interval to obtain the estimation of the rank(x).

483 Based on the observations, DCS and DCM quantile sketches keeps $\log U$ number of linear sketches,
484 one for each dyadic interval. As a result, to update an item x with value $v \in \{-1, +1\}$, DCS and
485 DCM need to update $\log U$ levels: they first map item x to a dyadic interval for the level and then
486 update the corresponding linear sketch, as shown in Algorithm 4. To estimate the rank of an item,
487 DCS and DCM first break the range into at most $\log U$ dyadic intervals and then query the frequency
488 for each interval from the corresponding linear sketch, as shown in Algorithm 5.