
Private Federated Frequency Estimation: Adapting to the Hardness of the Instance

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In *federated frequency estimation* (FFE), multiple clients work together to estimate
2 the frequency of their local data by communicating with a server, while respecting
3 the security constraints of Secure Summation (SecSum) where the server can only
4 access the sum of client-held vectors. For FFE with a single communication
5 round, it is known that *count sketch* is nearly information-theoretically optimal
6 [8]. However, when multiple communication rounds are allowed, we propose a
7 new sketching algorithm that is *provably* more accurate than a naive adaptation
8 of count sketch. Furthermore, we show that both our sketch algorithm and count
9 sketch can achieve better accuracy when the problem instance is simpler. Therefore,
10 we propose a two-phase approach to enable the use of a smaller sketch size for
11 simpler problems. Finally, we provide mechanisms to make our proposed algorithm
12 differentially private. We verify the superior performance of our methods through
13 experiments conducted on several largescale datasets.

14 1 Introduction

15 In many distributed learning applications, a server seeks to compute population information about
16 data that is distributed across multiple clients (users). For example, consider a distributed frequency
17 estimation problem where there are n clients, each holding a local data from a domain of size d ,
18 and a server that aims to estimate the frequency of the items from the n clients with a minimum
19 communication cost. This task can be done efficiently by letting each client *binary encode* their data
20 and send the encoding to the server, at a local communication bandwidth cost of $\log(d)$ bits. Provided
21 with the binary encoding, the server can faithfully decode *each* local data and compute the global
22 frequency vector (i.e., the normalized histogram vector).

23 However, the local data could be sensitive or private, and the clients may wish to keep it hidden
24 from the server. The above binary encoding communication method, unfortunately, allows the server
25 to observe each individual local data, and therefore may not satisfy the users' privacy concerns.
26 *Federated Analytics* (FA) [16, 18] addresses this issue by developing new methods that enable the
27 server to learn population information about the clients while preventing the server from prying on
28 any individual local data. In particular, a cryptographic multi-party computation protocol, *Secure*
29 *Summation* (SecSum) [1], has become a widely adopted solution to provide data minimization
30 guarantees for FA [3]. Specifically, SecSum sets up a communication protocol between clients and
31 the server, which injects carefully designed additive noise to each data that cancels out when *all of*
32 *the local data is summed together*, but blurs out (information theoretically) each individual local data
33 otherwise. Under SecSum, the server is able to faithfully obtain the correct summation of the data from
34 all clients but is unable to read a single local data. *Federated frequency estimation* (FFE) problems
35 refer to the distributed frequency estimation problems under the constraint of SecSum. Clearly, the
36 binary encoding method is not compatible with SecSum, because when the binary encoding is passed

37 to the server through SecSum, the server only gets the summation of the binary encodings of the
38 users’ data, which does not provide sufficient information for computing the global frequency vector.

39 A naive approach to FFE is by employing *one-hot encoding*: each client encodes its local data into
40 a d -dimensional one-hot vector that represents the local frequency vector and sends it to the server
41 through SecSum. Then the server observes the summation of the local frequency vectors using
42 SecSum and scales it by the number of clients to obtain the true frequency vector. However, the
43 one-hot encoding approach costs $\Theta(d \log(n))$ bits of communication bandwidth. This is because
44 SecSum adds noise from a field of size $\Theta(n)$ to each component of the d -dimensional local frequency
45 vector [1]. With a linear dependence on domain size d , the one-hot encoding approach is inefficient
46 for large domain problems, especially when the domain size exceeds the number of clients ($d > n$).
47 In what follows, we will focus on this regime and assume that $d > n$.

48 Recently, linear compression methods were applied to mitigate the high communication cost issue
49 for FFE with large domains [7, 8]. The idea is to first *linearly compress* the local frequency
50 vector into a lower dimensional vector before sending it to the server through SecSum; as linear
51 compression operators commute with the summation operator, the server equivalently observes a
52 linearly compressed global frequency vector through SecSum (after rescaling by the number of clients).
53 The server then applies standard decoding methods to approximately recover the global frequency
54 vector from the linearly compressed one. In particular, Chen et al. [8] show that CountSketch [6]
55 (among other sparse recover methods) can be used as a linear compressor for the above purpose,
56 which leads to a communication bandwidth cost of $\mathcal{O}(n \log(d) \log(n))$ bits. Therefore when $d >$
57 n , CountSketch achieves a saving in local communication bandwidth compared to the one-hot
58 encoding method that requires $\Theta(d \log(n))$ bits. Moreover, Chen et al. [8] show that for FFE with a
59 single communication round, an $\Omega(n \log(d))$ local communication cost is information-theoretically
60 *unavoidable* for worst-case data distributions, i.e., without making additional assumptions on the
61 global frequency vector.

62 **Contributions.** In this work, we make three notable extensions to CountSketch for FFE problems.

- 63 1. Firstly, we show that the way Chen et al. [8] set up the sketch size (linear in the number of clients
64 n) is often *pessimistic* (see Corollary 2.4). In fact, in the streaming literature, the estimation error
65 afforded by CountSketch is known to adapt to the tail norm of the global frequency vector [13],
66 which is often sub-linear in n . Motivated by this, we provide an easy-to-use, two-phase approach
67 that allows practitioners to determine the necessary sketch size by automatically adapting to the
68 hardness the FFE problem instance.
- 69 2. Secondly, we consider FFE with multiple communication rounds, which better models practical
70 deployments of FA where aggregating over (hundreds of) millions of clients in a single round is not
71 possible due to device availability and limited server bandwidth. We propose a new multi-round
72 sketch algorithm called HybridSketch that *provably* performs better than simple adaptations of
73 CountSketch in the multi-round setting, leading to further improvements in the communication
74 cost. Quite surprisingly, we show that HybridSketch adapts to the tail norm of a *heterogeneity*
75 *vector* (see Theorem 3.2). Moreover, the tail of the heterogeneity vector is always no heavier,
76 and could be much lighter, than that of the global frequency vector, explaining the advantage
77 of HybridSketch. For instance, on the C4 dataset [4] with a domain size of $d = 150, 868$ and
78 $150, 000$ users, we show that our method can reduce the sketch size by 83% relative to simple
79 sketching methods.
- 80 3. Finally, we extend the Gaussian mechanism for CountSketch proposed by Pagh and Thorup
81 [14], Zhao et al. [17] to the multi-round FFE setting to show how our sketching methods can be
82 made differentially private [10]. We also characterize the trade-offs between accuracy and privacy
83 for our proposed method.

84 We conclude by verifying the performance of our methods through experiments conducted on several
85 large-scale datasets. All proofs and additional experimental results are deferred to the appendices.

86 2 Adapting CountSketch to the Hardness of the Instance

87 In this part, we focus on single-round FFE and show how CountSketch can achieve better results
88 when the underlying problem is simpler. Motivated by this, we also provide a two-phase method for

89 auto-tuning the hyperparameters of CountSketch, allowing it to automatically adapt to the hardness
 90 of the instance.

91 **Single-Round FFE.** Consider n clients, each holding an item from a discrete domain of size d .
 92 The items are denoted by $x_t \in [d]$ for $t = 1, \dots, n$. Then the frequency of item j is denoted by

$$f_j := \frac{1}{n} \sum_{t=1}^n \mathbb{1}[x_t = j].$$

93 We use \mathbf{x}_t to denote the one-hot representation of x_t , i.e., $\mathbf{x}_t = \mathbf{e}_{x_t}$ where $(\mathbf{e}_t)_{t=1}^d$ refers to the
 94 canonical basis. Then the frequency vector can be denoted by

$$\mathbf{f} := (f_1, \dots, f_d)^\top = \frac{1}{n} \sum_{t=1}^n \mathbf{x}_t \in [0, 1]^d.$$

95 In single-round FFE, the n clients communicate with a server once under the constraint of SecSum,
 96 and aim to estimate the frequency vector \mathbf{f} . Note that SecSum ensures that the server can only observe
 97 the sum of the local data.

98 **Count Sketch.** CountSketch is a classic streaming algorithm that dates back to [6]. In the
 99 literature of streaming algorithms, CountSketch has been extensively studied and is known to be
 100 able to adapt to the hardness of the problem instance. Specifically, CountSketch of a fixed size
 101 induces an estimation error adapting to the tail norm of the global frequency vector [13].

102 A recent work by Chen et al. [8] apply CountSketch to single-round FFE. See Algorithm 2 in
 103 Appendix A for details. They show that CountSketch approximately solves single-round FFE
 104 with a communication cost of $\mathcal{O}(n \log(d) \log(n))$ bits per client. Moreover, they show $\Omega(n \log(d))$
 105 bits of communication per client is unavoidable for worst-case data distributions (unless additional
 106 assumptions are made), confirming its near optimality. However, the results by Chen et al. [8]
 107 are *pessimistic* as they ignore the ability of CountSketch to adapt to the hardness of the problem
 108 instance. In what follows, we show how the performance of CountSketch can be improved when
 109 the underlying problem becomes simpler.

110 We first present a problem-dependent accuracy guarantee for CountSketch of a fixed size, $L \times W$,
 111 that gives the sharpest bound to our knowledge. The bound is due to Minton and Price [13] and is
 112 restated for our purpose.

113 **Proposition 2.1** (Restated Theorem 4.1 in Minton and Price [13]). *Let $(\hat{f}_j)_{j=1}^d$ be estimates produced
 114 by CountSketch (see Algorithm 2). Then for each $p \in (0, 1)$, $W \geq 2$ and $L \geq \log(1/p)$, it holds
 115 that: for each $j \in [d]$, with probability at least $1 - p$,*

$$|\hat{f}_j - f_j| < C \cdot \sqrt{\frac{\log(1/p)}{L}} \cdot \frac{1}{W} \cdot \sum_{i>W} (f_i^*)^2,$$

116 where $(f_i^*)_{i \geq 1}$ refers to $(f_i)_{i \geq 1}$ sorted in non-increasing order, and $C > 0$ is an absolute constant.

117 For the concreteness of discussion, we will focus on ℓ_∞ as a measure of estimation error in the
 118 remainder of the paper. Our discussions can be easily extended to ℓ_2 or other types of error measures.
 119 Proposition 2.1 directly implies the following ℓ_∞ -error bounds for CountSketch (by an application
 120 of union bound).

121 **Corollary 2.2** (ℓ_∞ -error bounds for CountSketch). *Consider Algorithm 2. Then for each $p \in (0, 1)$,
 122 $L = \log(d/p)$ and $W \geq 2$, it holds that: with probability at least $1 - p$,*

$$\|\text{dec}(\cdot) - \mathbf{f}\|_\infty < C \cdot \sqrt{\frac{1}{W}} \cdot \sum_{i>W} (f_i^*)^2, \quad (1)$$

123 where $C > 0$ is an absolute constant. In particular, (1) implies that

$$\|\text{dec}(\cdot) - \mathbf{f}\|_\infty < \frac{C}{W}.$$

124 According to Corollary 2.2, the estimation error will be smaller when the underlying frequency vector
 125 $(f_i^*)_{i \geq 1}$ has a lighter tail. Said differently, CountSketch needs less communication bandwidth when
 126 the global frequency vector has a lighter tail. Our next Corollary 2.3 precisely characterizes this
 127 adaptive property in terms of the required communication bandwidth. To show this, we will need the
 128 following definition on the *probable approximate correctness* of an estimate.

129 **Definition 1** ((τ, p) -correctness). An estimate $\hat{\mathbf{f}} := (\hat{f}_i)_{i=1}^d$ of the global frequency vector $\mathbf{f} :=$
 130 $(f_i)_{i=1}^d$ is (τ, p) -correct if

$$\mathbb{P}\left\{\|\hat{\mathbf{f}} - \mathbf{f}\|_\infty := \max_i |\hat{f}_i - f_i| > \tau\right\} < p.$$

131 **Corollary 2.3** (Oracle sketch size). Fix parameters $\tau, p \in (0, 1)$. Then for CountSketch (see
 132 Algorithm 2) to produce an (τ, p) -correct estimate, it suffices to set the sketch size to $L = \log(d/p)$
 133 and

$$W = C \cdot \min \left\{ \left(\#\{f_i : f_i \geq \tau\} + \frac{1}{\tau^2} \cdot \sum_{f_i < \tau} f_i^2 \right), n \right\}, \quad (2)$$

134 where $C > 0$ is an absolute constant. In particular, the width W in (2) satisfies

$$W \leq W_{\text{worst}} := C \cdot \min \{2/\tau, n\}. \quad (3)$$

135 Corollary 2.3 suggests that the sketch size can be made smaller if the underlying frequency vector
 136 has a lighter tail. When translated to the communication bits per client (that is $\mathcal{O}(L \cdot W \cdot \log(n))$),
 137 where $\log(n)$ accounts for the cost of SecSum), Corollary 2.3 implies that CountSketch requires

$$\mathcal{O}\left(\min \left\{ \#\{f_i \geq \tau\} + \frac{1}{\tau^2} \sum_{f_i < \tau} f_i^2, n \right\} \log(d) \log(n)\right) \leq \mathcal{O}(\min\{1/\tau, n\} \log(d) \log(n)) \quad (4)$$

138 bits of communication per client to be (τ, p) -correct. In the worst case where $(f_i)_{i=1}^d$ is $\Theta(n)$ -sparse
 139 and $\tau = \mathcal{O}(1/n)$, (4) nearly matches the $\Omega(n \log(d))$ information-theoretic worst-case communica-
 140 tion cost shown in Chen et al. [8], ignoring the $\log(n)$ factor from SecSum. However, in practice,
 141 $(f_i)_{i=1}^d$ has a fast-decaying tail, and (4) suggests that CountSketch can use less communication for
 142 solving the problem. We provide the following examples for a better illustration of the sharp contrast
 143 between the worst and typical cases.

144 **Corollary 2.4** (Examples). Fix parameters $\tau, p \in (0, 1)$. Consider Algorithm 2 with sketch length
 145 $L = \log(d/p)$. Then in each case for Algorithm 2 to produce an (τ, p) -correct estimate for $\tau > 1/n$:

- 146 1. When $f_i \propto 2^{-i}$, it suffices to set $W = \Theta(\log(1/\tau))$.
- 147 2. When $f_i \propto i^{-a}$ for $a > 1$, it suffices to set $W = \Theta(\tau^{-1/a})$.
- 148 3. When $f_i \propto i^{-1} \log^{-b}(i)$ for $b > 1$, it suffices to set $W = \Theta(\tau^{-1} \log^{-b}(1/\tau))$.
- 149 4. When $f_i = 10/n$ for $i = 1, \dots, n/10$, it suffices to set $W = \Theta(1/\tau)$.

150 **A Two-Phase Method for Hyperparameter Setup.** Corollary 2.3 allows to use CountSketch
 151 with a smaller width for an easier single-round FFE problem, saving communication bandwidth.
 152 However, the sketch size formula given by (2) in Corollary 2.3 relies on crucial information of the
 153 frequency $(f_i)_{i \geq 1}$, i.e., $\#\{f_i : f_i \geq \tau\}$ and $\sum_{f_i < \tau} f_i^2$, which are unknown to the engineer who sets
 154 the sketch size. Thus, it is unclear if and how these gains can be realized in practical deployments.

155 We resolve this quandary by observing that in practice, the frequency vector often follows Zipf's
 156 law [5, 15]. This motivates us to conservatively model the global frequency vector by a polynomial
 157 with unknown parameters. By doing so, we can first run a small CountSketch to collect data from a
 158 (randomly sampled) fraction of the clients for estimating the parameters. Then based on the estimated
 159 parameter, we can set up an appropriate sketch size for a CountSketch to solve the FFE problem.
 160 This two-phase method is formally stated as follows.

161 We approximate the (sorted) global frequency vector $(f_i^*)_{i=1}^d$ by a polynomial [5] with two parameters
 162 $\alpha > 0$ and $\beta > 0$, such that

$$f_i^* \approx \text{poly}(i; \alpha, \beta), \quad \text{poly}(i; \alpha, \beta) := \begin{cases} \beta \cdot i^{-\alpha}, & i \leq i^*; \\ 0, & i > i^*, \end{cases}$$

163 where $i^* := \max\{i : \sum_{j=1}^i \beta \cdot j^{-\alpha} \leq 1\}$ is set such that $\text{poly}(i; \alpha, \beta)$ is a valid frequency vector.
 164 Here's an executive summary of the proposed approach for setting the sketch size.

- 165 1. Randomly select a subset of clients (e.g., 5,000 out of 10^6 .)
- 166 2. Fix a small sketch (e.g., 16×100) and run Algorithm 2 with the subset of clients to obtain an
 167 estimate (\tilde{f}_i) .
- 168 3. Use the top- k values (e.g., top 20) from \tilde{f}_i to fit a polynomial with parameter α and β (under
 169 squared error).
- 170 4. Solve Equation (4) under the approximation that $f_i^* \approx \beta \cdot i^{-\alpha}$ and output W according to the result.

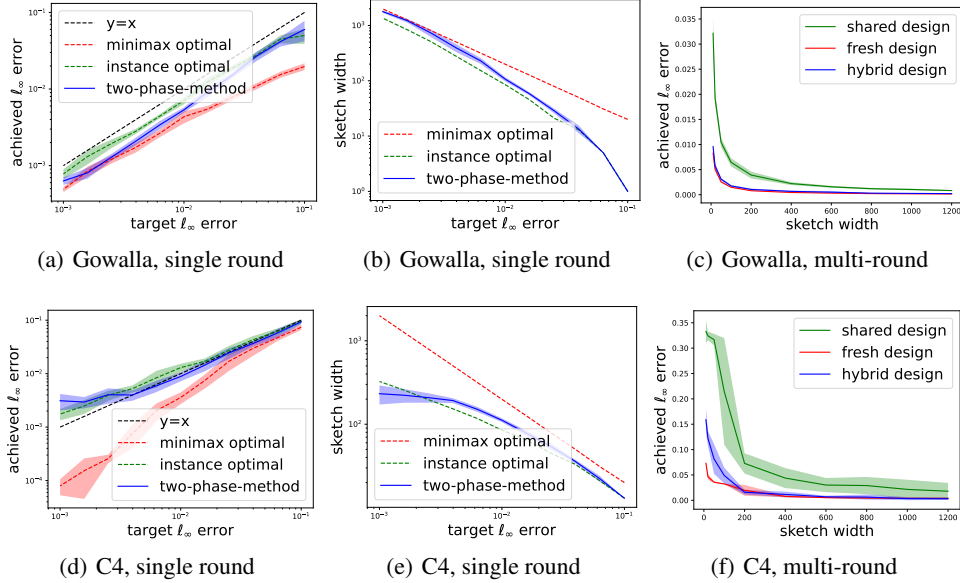


Figure 1: Single-round and multi-round FFE simulations. Subfigures (a) and (b) compare different hyperparameter strategies for CountSketch in a single-round FFE problem on the Gowalla dataset [9]. Subfigure (c) compares three sketch methods in a multi-round FFE problem on the Gowalla dataset. Subfigures (d), (e), and (f) are counterparts of subfigures (a), (b), and (c), respectively, but on the C4 [4] dataset.

171 **Experiments.** We conduct two sets of experiments to verify our methods. In the first set of
 172 experiments, we simulate a single-round FFE problem with the Gowalla dataset [9]. The dataset
 173 contains 6,442,892 lists of location information. We first construct a domain of size $d = 175,000$,
 174 which corresponds to a grid over the US map. Then we sample $n = d/10 = 17,500$ lists of
 175 the location information (that all belong to the domain created) to represent the data of n clients,
 176 uniformly at random. This way, we set up a single-round FFE problem with $n = 17,500$ clients in
 177 a domain of size $d = 175,000$. In the experiments, we fix the confidence parameter to be $p = 0.1$
 178 and the sketch length to be $L = \ln(2d/p) \approx 16$. The targeted ℓ_∞ -error τ is chosen evenly from
 179 $(10^{-3}, 10^{-1})$. We only test $\tau > 20/n$ because it is less important to estimate frequencies over
 180 items with small counts (say, 20). For CountSketch, we compute sketch width with three strategies,
 181 using (2) (called “instance optimal”), using (3) (called “minimax optimal”), and using the two-phase
 182 method. We set all constant factors to be 2. The results are presented in Figures 1(a) and (b). We
 183 observe that the “minimax optimal” way of hyperparameter choice is in fact suboptimal in practice,
 184 and is improved by the “instance optimal” and the two-phase strategies.

185 In the second set of experiments, we run simulations on the “Colossal Clean Crawled Corpus” (C4)
 186 dataset [4], which consists of clean English text scraped from the web. We treat each domain in the
 187 dataset as a user and calculate the number of examples each user has. The domain size $d = 150,868$,
 188 which is the maximum example count per user. We randomly sample $n = 150,000$ users from the
 189 dataset. We fix the confidence parameter to be $p = 0.1$ and the sketch length to be $L = 5$. Other
 190 parameters are the same as the Gowalla dataset. The results are presented in Figures 1(d) and (e), and
 191 are consistent with what we have observed in the Gowalla simulations.

192 3 Sketch Methods for Multi-Round Federated Frequency Estimation

193 In practice, having all clients participate in a single communication round is impractical due to the
 194 large number of devices, their unpredictable availability, and limited server bandwidth [2]. This
 195 motivates us to consider a multi-round FFE setting.

196 **Multi-Round FFE.** Consider a FFE problem with M rounds of communication. In each round,
 197 n clients participate, each holding an item from a universe of size d . The items are denoted by
 198 $x_t^{(m)} \in [d]$, where $t \in [n]$ denotes the client index and $m \in [M]$ denotes the round index. For

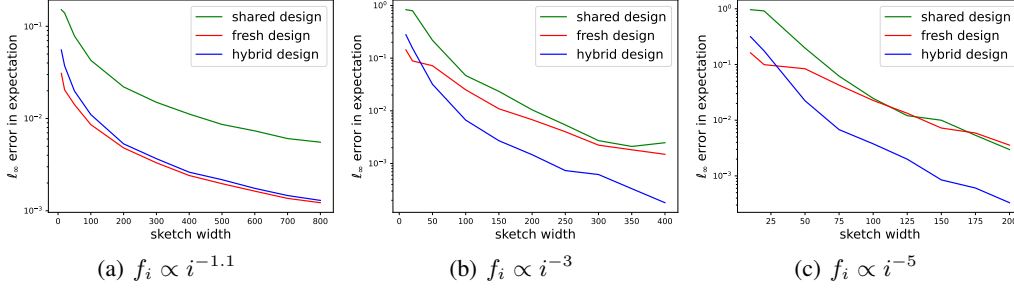


Figure 2: Shared vs. Hybrid vs. Fresh Sketches. We refer the reader to Section 3 for the definitions of the three methods. We compute the expected ℓ_∞ -error for shared/hybrid/fresh sketches for a homogeneous, multi-round FFE problem. The domain size is $d = 10^5$. The number of rounds is $M = 10$. In all setups, the sketch length is fixed to $L = 5$. In every setting, the ℓ_∞ error is averaged with 1,000 random repeats for simulating the expectation. In the case when the global frequency vector is a low-degree polynomial, hybrid sketch performs similarly to fresh sketch, and both are better than shared sketch. As long as the global frequency vector is a slightly higher degree polynomial (e.g., with a degree higher than 3), then hybrid sketch is significantly better than both shared and fresh sketches.

199 simplicity, we assume in each round a new set of clients participate. So in total there are $N = Mn$
 200 clients. Then the frequency of item j is now denoted by

$$f_j := \frac{1}{Mn} \sum_{m=1}^M \sum_{t=1}^n \mathbb{1} [x_t^{(m)} = j].$$

201 For the m -th round, the local frequency is denoted by $f_j^{(m)} := \frac{1}{n} \sum_{t=1}^n \mathbb{1} [x_t^{(m)} = j]$. Clearly, we
 202 have $f_j = \frac{1}{M} \sum_{m=1}^M f_j^{(m)}$. Similarly, we use $\mathbf{x}_t^{(m)}$ to denote the one-hot representation of $x_t^{(m)}$, i.e.,
 203 $\mathbf{x}_t^{(m)} = \mathbf{e}_{x_t^{(m)}}$ where $(\mathbf{e}_t)_{t=1}^d$ refers to the canonical basis. Then the frequency vector can be denoted
 204 by $\mathbf{f} := (f_1, \dots, f_d)^\top$. The aim is to estimate the frequency vector \mathbf{f} in a manner that is compatible
 205 with SecSum.

206 **Baseline Method 1: Shared Sketch.** A multi-round FFE problem can be reduced to a single-round
 207 FFE problem with a large communication. Specifically, one can apply the CountSketch with the
 208 same randomness for every round; after collecting all the sketches from the M round, one simply
 209 averages them. Due to the linearity of the sketching compress method, this is equivalent to a single
 210 round setting with $N = Mn$ clients. We refer to this method as *count sketch with shared hash design*
 211 (SharedSketch).

212 Thanks to the reduction idea, we can obtain the error and sketch size bounds for SharedSketch by
 213 applying Corollaries 2.2 and 2.3 to SharedSketch by replacing n by $N = Mn$,

214 **Baseline Method 2: Fresh Sketch.** A multi-round FFE problem can also be broken down to M
 215 independent single-round FFE problems. Specifically, one can apply *independent* CountSketch in
 216 each round, and decode M local estimators for the M local frequency vectors. As the CountSketch
 217 produces an unbiased estimator, one can show that the average of the M local estimators is an
 218 unbiased estimator for the global frequency vector. We call this method *count sketch with fresh hash*
 219 *design* (FreshSketch). We provide the following bound for FreshSketch. The proof of which is
 220 motivated by Huang et al. [12].

221 **Theorem 3.1** (Instance-specific bound for FreshSketch). *Let $(\hat{f}_j)_{j=1}^d$ be estimates produced by*
 222 *FreshSketch. Then for each $p \in (0, 1)$, $W \geq 1$ and $L \geq \log(1/p)$, it holds that: for each $j \in [d]$,*
 223 *with probability at least $1 - p$,*

$$|\hat{f}_j - f_j| < C \cdot \sqrt{\frac{\log(1/p) \log(M/p)}{L}} \cdot \frac{1}{W} \cdot \sum_{i>W} (F_i^*)^2,$$

224 where C is an absolute constant, and $(F_i^*)_{i=1}^d$ are defined as in Theorem 3.2.

Algorithm 1 HYBRID SKETCH FOR FEDERATED FREQUENCY ESTIMATION

Require: The number of rounds M . $N = Mn$ clients with local data $x_t^{(m)} \in [d]$ for $m \in [M]$ and $t \in [n]$. Sketch length L and width W .

1: The server prepares independent hash functions and broadcasts them to each client:

$$h_\ell : [d] \rightarrow [W], \sigma_\ell^{(m)} : [d] \rightarrow \{\pm 1\} \text{ for } \ell \in [L], m \in [M].$$

2: **for** Round $m = 1, \dots, M$ in parallel **do**

3: **for** Client $t = 1, \dots, n$ in parallel **do**

4: Client (m, t) encodes the local data $x_t^{(m)}$ to $\text{enc}^{(m)}(x_t^{(m)}) \in \mathbb{R}^{L \times W}$ where

$$\left(\text{enc}^{(m)}(x_t^{(m)})\right)_{\ell, k} = \mathbb{1} \left[h_\ell(x_t^{(m)}) = k \right] \cdot \sigma_\ell^{(m)}(x_t^{(m)}) \text{ for } \ell \in [L], k \in [W].$$

5: Client (m, t) sends $\text{enc}^{(m)}(x_t^{(m)})$ to SecSum.

6: **end for**

7: SecSum receives $\left(\text{enc}^{(m)}(x_t^{(m)})\right)_{t=1}^n$ and reveals the sum $\sum_{t=1}^n \text{enc}^{(m)}(x_t^{(m)})$ to the server.

8: **end for**

9: **for** Item $j = 1, \dots, d$ in parallel **do**

10: Server produces $M \times L$ estimators for f_j :

$$\text{dec}(j; m, l) := \sigma_\ell^{(m)}(j) \cdot \left(\frac{1}{n} \sum_{t=1}^n \text{enc}^{(m)}(x_t^{(m)}) \right)_{\ell, h_\ell(j)} \text{ for } m \in [M], \ell \in [L].$$

11: Server computes the median over $\ell \in [L]$ of the averages over $m \in [M]$ of the estimators:

$$\text{dec}(j) := \text{median} \left\{ \frac{1}{M} \sum_{m=1}^M \text{dec}(j; m, l), \ell \in [L] \right\}.$$

12: **end for**

13: **return** $(\text{dec}(j))_{j=1}^d$ as estimate to $(f_j)_{j=1}^d$.

225 **Hybrid Sketch.** Both SharedSketch and FreshSketch are reducing a multi-round FFE problem
 226 into single-round FFE problem(s). Instead, we show a more comprehensive sketching method, called
 227 *count sketch with hybrid hash design* (HybridSketch), that solves a multi-round FFE problem as
 228 a whole. HybridSketch is presented as Algorithm 1. Specifically, HybridSketch generates M
 229 sketches that share a set of bucket hashes but use independent sets of sign hashes. Then in the
 230 m -th communication round, participating clients and the server communicate by the CountSketch
 231 algorithm based on the m -th sketch, so the server observes the summation of the sketched data
 232 through SecSum. After collecting M summations of the sketched local data, the server first computes
 233 averages over different rounds for *variance reduction*, then computes the median over different repeats
 234 (or sketch rows) for *success probability amplification*. We provide the following problem-dependent
 235 bound for HybridSketch.

236 **Theorem 3.2** (Instance-specific bound for HybridSketch). Let $(\hat{f}_j)_{j=1}^d$ be estimates produced by
 237 HybridSketch (see Algorithm 1). Define a heterogeneity vector $(F_i)_{i=1}^d$ by

$$F_i := \frac{1}{M} \sqrt{\sum_{m=1}^M (f_i^{(m)})^2}, \quad i = 1, \dots, d.$$

238 Clearly, it holds that $F_i \leq f_i$ for every $i \in [d]$. Let $(F_i^*)_{i \geq 1}$ be $(F_i)_{i \geq 1}$ sorted in non-increasing
 239 order. Then for each $p \in (0, 1)$, $W \geq 1$ and $L \geq \log(1/p)$, it holds that: for each $j \in [d]$, with
 240 probability at least $1 - p$,

$$|\hat{f}_j - f_j| < C \cdot \sqrt{\frac{\log(1/p)}{L}} \cdot \frac{1}{W} \cdot \sum_{i > W} (F_i^*)^2,$$

241 where C is an absolute constant.

242 **Hybrid Sketch vs. Fresh Sketch.** By comparing Theorem 3.2 with Theorem 3.1, we see that, with
 243 the same sketch size, the estimation error of HybridSketch is smaller than that of FreshSketch
 244 by a factor of $\sqrt{\log(M/p)}$. This provides theoretical insights that HybridSketch is superior to
 245 FreshSketch in terms of adapting to the instance hardness in multi-round FFE settings. This is also
 246 verified empirically by Figure 2.

247 **Hybrid Sketch vs. Shared Sketch.** We now
 248 compare the performance of HybridSketch
 249 and SharedSketch by comparing Theorem 3.2
 250 and Proposition 2.1 (under a revision of replac-
 251 ing n with $N = Mn$). Note that

$$F_i = \frac{1}{M} \sqrt{\sum_{m=1}^M (f_i^{(m)})^2} \leq \frac{1}{M} \sum_{m=1}^M f_i^{(m)} = f_i.$$

252 So with the same sketch size, HybridSketch
 253 achieves an error that is no worse than csc
 254 in every case. Moreover, in the *homogeneous*
 255 case where all local frequency vectors are equivalent
 256 to the global frequency vector, i.e., $\mathbf{f}^{(m)} \equiv \mathbf{f}$ for
 257 all m , then it holds that $F_i = f_i/\sqrt{M}$. So in the
 258 homogeneous case, HybridSketch achieves an
 259 error that is smaller than that of csc by a factor
 260 of $1/\sqrt{M}$. In the general cases, the local fre-
 261 quency vectors are not perfectly homogeneous,
 262 then the improvement of HybridSketch over
 263 SharedSketch will depend on the *heterogeneity*
 264 of these local frequency vectors.

265 **Experiments.** We conduct three sets of experiments to verify our understandings about these
 266 sketches methods for multi-round FFE.

267 In the first sets of experiments, we simulate a multi-round FFE problem in homogeneous settings,
 268 where in every round the local frequency vectors are exactly the same. More specially, we set a
 269 domain size $d = 10^5$, a number of rounds $M = 10$ and test three different cases, where all the
 270 local frequency vectors are the same and (hence also the global frequency vector) are proportional
 271 to $(i^{-1.1})_{i=1}^d$, $(i^{-2})_{i=1}^d$ and $(i^{-5})_{i=1}^d$, respectively. In all the settings, we fix the sketch length to
 272 $L = 5$. In each experiment, we measure the expected ℓ_∞ -error of each method with the averaging
 273 over 1,000 independent repeats. The results are plotted in Figure 2. We can observe that: for
 274 low-degree polynomials, HybridSketch is nearly as good as FreshSketch and both are better
 275 than SharedSketch. But for slightly high degree polynomials (with a degree of 3), HybridSketch
 276 already outperforms both FreshSketch and SharedSketch. The numerical results are consistent
 277 with our theoretical analysis.

278 In the second sets of experiments, we simulate a multi-round FFE problem with the Gowalla dataset
 279 [9]. Similar to previously, we construct a domain of size $d = 175,000$, which corresponds to a grid
 280 over the US map. Then we sample $N = d = 175,000$ lists of the location information (that all
 281 belong to the domain created) to represent the data of N clients, uniformly at random. We set the
 282 number of rounds to be $M = 10$. In each round, $n = N/M = 17,500$ clients participate. The results
 283 are presented in Figure 1(c). Here, the frequency and heterogeneity vectors have heavy tails, so
 284 HybridSketch and FreshSketch perform similarly and both are better than SharedSketch. This
 285 is consistent with our theoretical understanding.

286 In the third sets of experiments, we run simulations on the C4 [4] dataset. Similar to the single
 287 round simulation, the domain size $d = 150,868$. We randomly sample $N = 150,000$ users from the
 288 dataset. The number of rounds $M = 10$, and in each round, $n = N/10 = 15,000$ clients participate.
 289 The results are provided in Figures 1(f) and 3. Here, the frequency and heterogeneity vectors have
 290 moderately light tails, and Figure 3 already suggests that HybridSketch produces an estimate that
 291 has a better shape than that produced by FreshSketch and SharedSketch, verifying the advantages
 292 of HybridSketch.

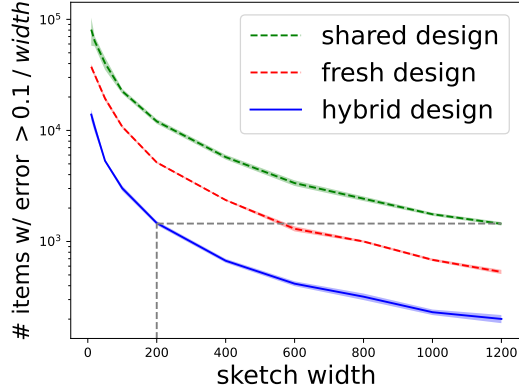


Figure 3: The number of items with error greater than $0.1/\text{width}$ for Shared, Hybrid, and Fresh Sketches with C4 dataset. HybridSketch with a width of 200 achieves roughly the same error as SharedSketch with a width of 1200 and Fresh sketch with a width of 600.

293 4 Differentially Private Sketches

294 While SecSum provides security guarantees, it does not provide differential privacy guarantees. In
 295 this part, we discuss a simple modifications to the sketching algorithms to make them provably
 296 differentially private.

297 **Definition 2** ((ϵ, δ) -DP [10]). Let $\text{alg}(\cdot)$ be a randomized algorithm that takes a dataset \mathcal{D} as its
 298 input. Let \mathbb{P} be its probability measure. $\text{alg}(\cdot)$ is (ϵ, δ) -DP if: for every pair of neighboring datasets
 299 \mathcal{D} and \mathcal{D}' , it holds that

$$\mathbb{P}\{\text{alg}(\mathcal{D}) \in \mathcal{E}\} < e^\epsilon \cdot \mathbb{P}\{\text{alg}(\mathcal{D}') \in \mathcal{E}\} + \delta.$$

300 In our case, a dataset corresponds to all participated clients (or their data), and two neighboring
 301 datasets should be regarded as two sets of clients (local data) that only differ in a single client (local
 302 data). The algorithm refers to all procedures before releasing the final frequency estimate, and all the
 303 intermediate computation is considered private and is not released.

304 We focus on HybridSketch as a representative algorithm. The DP mechanism can also be extended to
 305 the other sketching algorithms. Specifically, we use a DP mechanism that adds independent Gaussian
 306 noise to each entry of the sketching matrix, which is initially proposed for making CountSketch
 307 differentially private by Pagh and Thorup [14], Zhao et al. [17].

308 We provide the following theorem characterizing the trade-off between privacy and accuracy.

309 **Theorem 4.1** (DP-hybrid sketch). *Consider a modified Algorithm 1, where we add to each entry of
 310 the sketching matrix an independent Gaussian noise, $\mathcal{N}(0, c_0 \cdot \sqrt{L \log(1/\delta)}/\epsilon)$, where $c_0 > 0$ is
 311 a known constant. Suppose that $L = \log(d/p)$ and $W \geq 2$. Then the final output of the modified
 312 Algorithm 1, denoted by $(\hat{f}_j)_{j=1}^d$, is (ϵ, δ) -DP for $\epsilon < 1$ and $\delta < 0.1$. Moreover, with probability at
 313 most $1 - p$, it holds that*

$$\max_j |\hat{f}_j - f_j| < C \cdot \left(\sqrt{\frac{\sum_{i>W} (F_i^*)^2}{W}} + \frac{\sqrt{\log(d/p) \log(1/\delta)}}{n\sqrt{M}\epsilon} \right),$$

314 where $C > 0$ is an absolute constant and $(F_i^*)_{i=1}^d$ are as defined in Theorem 3.2.

315 It is worth noting that if the number of clients per round (n) is fixed, then a larger number of rounds
 316 M improves both the estimation error and the DP error. However if the total number of clients
 317 ($N = Mn$) is fixed, then a larger number of rounds M improves the estimation error but makes the
 318 DP error worse.

319 When $M = 1$, Theorem 4.1 recovers the bounds for differentially private CountSketch in
 320 Pagh and Thorup [14], Zhao et al. [17] and Theorem 5.1 in Chen et al. [8]. Moreover, Chen
 321 et al. [8] shows that in single-round FFE, for any algorithm that achieves an ℓ_∞ -error smaller
 322 than $\tau := \mathcal{O}(\sqrt{\log(d) \log(1/\delta)}/(n\epsilon))$, in the worse case, each client must communicate $\Omega(n \cdot$
 323 $\min\{\sqrt{\log(d)/\log(1/\delta)}, \log(d)\})$ bits (see Their Corollary 5.1). In comparison, According to Theo-
 324 rem 4.1 and Corollary 2.3, the differentially private CountSketch can achieve an ℓ_∞ -error smaller
 325 than τ with length $L \approx \log(d)$ and width

$$W = C \cdot \min \left\{ \left(\#\{f_i : f_i \geq \tau\} + \frac{1}{\tau^2} \cdot \sum_{f_i < \tau} f_i^2 \right), n \right\} \leq C \cdot \min\{2/\tau, n\},$$

326 resulting in a per-client communication of $\mathcal{O}(WL \log(n))$ bits, which matches the minimax lower
 327 bound in Chen et al. [8] ignoring a $\log(n)$ factor, but could be much smaller in non-worst cases where
 328 $(f_i)_{i=1}^d$ decays fast.

329 5 Conclusion

330 We make several novel extensions to the count sketch method for federated frequency estimation
 331 with one or more communication rounds. In the single round setting, we show that count sketch
 332 can achieve better communication efficiency when the underlying problem is simpler. We provide a
 333 two-phase approach to automatically select a sketch size that adapts to the hardness of the problem. In
 334 the multiple rounds setting, we show a new sketching method that provably achieves better accuracy
 335 than simple adaptations of count sketch. Finally, we adapt the Gaussian mechanism to make the hybrid
 336 sketching method differentially private.

337 References

- 338 [1] Kallista Bonawitz, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H Brendan McMahan,
339 Sarvar Patel, Daniel Ramage, Aaron Segal, and Karn Seth. Practical secure aggregation for
340 federated learning on user-held data. *arXiv preprint arXiv:1611.04482*, 2016.
- 341 [2] Kallista Bonawitz, Hubert Eichner, Wolfgang Grieskamp, Dzmitry Huba, Alex Ingerman,
342 Vladimir Ivanov, Chloé Kiddon, Jakub Konečný, Stefano Mazzocchi, Brendan McMahan,
343 Timon Van Overveldt, David Petrou, Daniel Ramage, and Jason Roselander. Towards federated
344 learning at scale: System design. In *MLSys*. mlsys.org, 2019.
- 345 [3] Kallista Bonawitz, Peter Kairouz, Brendan McMahan, and Daniel Ramage. Federated learning
346 and privacy: Building privacy-preserving systems for machine learning and data science on
347 decentralized data. *Queue*, 19(5):87–114, 2021.
- 348 [4] Samuel R. Bowman, Gabriel Angeli, Siddharth Jain, Jared Kaplan, Prafulla Dhariwal, Saurabh
349 Neelakantan, Jonathon Shlens, and Dario Amodei. C4: Colossal clean crawled corpus. *arXiv*
350 *preprint arXiv:2005.14165*, 2020.
- 351 [5] Volkan Cevher. Learning with compressible priors. *Advances in Neural Information Processing*
352 *Systems*, 22, 2009.
- 353 [6] Moses Charikar, Kevin Chen, and Martin Farach-Colton. Finding frequent items in data streams.
354 In *International Colloquium on Automata, Languages, and Programming*, pages 693–703.
355 Springer, 2002.
- 356 [7] Wei-Ning Chen, Christopher A Choquette-Choo, and Peter Kairouz. Communication efficient
357 federated learning with secure aggregation and differential privacy. In *NeurIPS 2021 Workshop*
358 *Privacy in Machine Learning*, 2021.
- 359 [8] Wei-Ning Chen, Ayfer Özgür, Graham Cormode, and Akash Bharadwaj. The communication
360 cost of security and privacy in federated frequency estimation. *arXiv preprint arXiv:2211.10041*,
361 2022.
- 362 [9] Eunjoon Cho, Seth A Myers, and Jure Leskovec. Friendship and mobility: user movement
363 in location-based social networks. In *Proceedings of the 17th ACM SIGKDD international*
364 *conference on Knowledge discovery and data mining*, pages 1082–1090, 2011.
- 365 [10] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to
366 sensitivity in private data analysis. In *Theory of Cryptography: Third Theory of Cryptography*
367 *Conference, TCC 2006, New York, NY, USA, March 4-7, 2006. Proceedings 3*, pages 265–284.
368 Springer, 2006.
- 369 [11] Alec Go, Richa Bhayani, and Lei Huang. Twitter sentiment classification using distant supervi-
370 sion. *CS224N Project Report, Stanford*, 1(12), 2009.
- 371 [12] Ziyue Huang, Yuan Qiu, Ke Yi, and Graham Cormode. Frequency estimation under multiparty
372 differential privacy: One-shot and streaming. *arXiv preprint arXiv:2104.01808*, 2021.
- 373 [13] Gregory T Minton and Eric Price. Improved concentration bounds for count-sketch. In
374 *Proceedings of the twenty-fifth annual ACM-SIAM symposium on Discrete algorithms*, pages
375 669–686. SIAM, 2014.
- 376 [14] Rasmus Pagh and Mikkel Thorup. Improved utility analysis of private counts sketch. *arXiv*
377 *preprint arXiv:2205.08397*, 2022.
- 378 [15] David M. W. Powers. Applications and explanations of Zipf’s law. In *New Methods in*
379 *Language Processing and Computational Natural Language Learning*, 1998. URL <https://aclanthology.org/W98-1218>.
- 380
- 381 [16] Daniek Ramage and Stefano Mazzocchi. Federated analytics: Collaborative
382 data science without data collection. [https://ai.googleblog.com/2020/05/](https://ai.googleblog.com/2020/05/federated-analytics-collaborative-data.html)
383 [federated-analytics-collaborative-data.html](https://ai.googleblog.com/2020/05/federated-analytics-collaborative-data.html), 2020.

- 384 [17] Fuheng Zhao, Dan Qiao, Rachel Redberg, Divyakant Agrawal, Amr El Abbadi, and Yu-Xiang
385 Wang. Differentially private linear sketches: Efficient implementations and applications. *arXiv*
386 *preprint arXiv:2205.09873*, 2022.
- 387 [18] Wennan Zhu, Peter Kairouz, Brendan McMahan, Haicheng Sun, and Wei Li. Federated heavy
388 hitters discovery with differential privacy. In *International Conference on Artificial Intelligence*
389 *and Statistics*, pages 3837–3847. PMLR, 2020.

A Count Sketch for Federated Frequency Estimation

Algorithm 2 COUNT SKETCH FOR FEDERATED FREQUENCY ESTIMATION

Require: n clients with local data $x_t \in [d]$ for $t = 1, \dots, n$. Sketch length L and width W .

1: The server prepares independent hash functions and broadcasts them to each client:

$$h_\ell : [d] \rightarrow [W], \sigma_\ell : [d] \rightarrow \{\pm 1\}, \text{ for } \ell \in [L].$$

2: **for** Client $t = 1, \dots, n$ in parallel **do**

3: Client t encodes the local data $x_t \in [d]$ to $\text{enc}(x_t) \in \mathbb{R}^{L \times W}$ where

$$(\text{enc}(x_t))_{\ell, k} = \mathbb{1}[h_\ell(x_t) = k] \cdot \sigma_\ell(x_t) \text{ for } \ell \in [L], k \in [W].$$

4: Client t sends $\text{enc}(x_t) \in \mathbb{R}^{L \times W}$ to SecSum.

5: **end for**

6: SecSum receives $(\text{enc}(x_t))_{t=1}^n$ and reveals the summation $\sum_{t=1}^n \text{enc}(x_t)$ to the server.

7: **for** Item $j = 1, \dots, d$ in parallel **do**

8: Server produces L estimators for f_j :

$$\text{dec}(j; \ell) := \sigma_\ell(j) \cdot \left(\frac{1}{n} \sum_{t=1}^n \text{enc}(x_t) \right)_{\ell, h_\ell(j)} \text{ for } \ell \in [L].$$

9: Server computes the median of the L estimators:

$$\text{dec}(j) := \text{median}\{\text{dec}(j; \ell) : \ell \in [L]\}.$$

10: **end for**

11: **return** $(\text{dec}(j))_{j=1}^d$ as estimate to $(f_j)_{j=1}^d$.

391 **B Additional Experiments**

392 **Sentiment-140.** We also run additional simulations on a Twitter dataset Sentiment-140 [11]. The
 393 dataset contains $d = 739,972$ unique words from $N = 659,497$ users. We randomly sample one
 394 word from each user to construct our experiment dataset. The number of rounds $M = 10$, and in
 395 each round, $n = N/10 = 65,949$ clients participate. Results are provided in Figure 4.

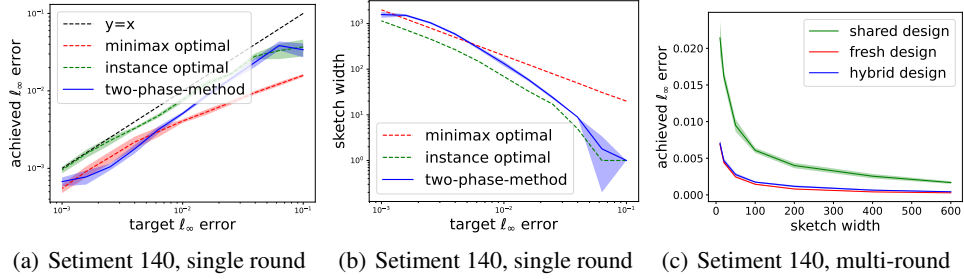


Figure 4: Single-round and multi-round FFE simulations on the Sentiment-140 dataset.

396 **Additional Plots for Single-Round FFE.** Figure 5 provides some additional results in our single-
 397 round FFE simulations.

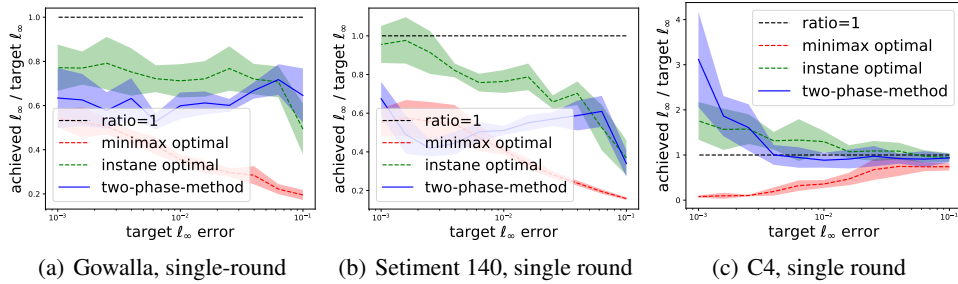


Figure 5: Single-round federated frequency estimation experiments.

398 **C Missing Proofs for Section 2**

399 **C.1 Proof of Proposition 2.1**

400 *Proof of Proposition 2.1.* We refer the reader to Theorem 4.1 in Minton and Price [13]. \square

401 **C.2 Proof of Corollary 2.2**

402 *Proof of Corollary 2.2.* From Proposition 1 we know that

$$\text{for every } j \in [d], \quad \mathbb{P}\left\{|\text{dec}(j) - f_j| > C \cdot \sqrt{\frac{\log(1/\delta)}{L} \cdot \frac{1}{W} \cdot \sum_{i>W} (f_i^*)^2}\right\} < \delta.$$

403 By union bound we have

$$\mathbb{P}\left\{\text{there exists } j \in [d], |\text{dec}(j) - f_j| > C \cdot \sqrt{\frac{\log(1/\delta)}{L} \cdot \frac{1}{W} \cdot \sum_{i>W} (f_i^*)^2}\right\} < d\delta.$$

404 Replacing δ with δ/d , setting $L = \log(d/\delta)$, and using the definition of ℓ_∞ -norm, we obtain

$$\mathbb{P}\left\{\|\text{dec}(\cdot) - \mathbf{f}\|_\infty > C \cdot \sqrt{\frac{1}{W} \cdot \sum_{i>W} (f_i^*)^2}\right\} < \delta.$$

405 We next show that:

$$\sqrt{\frac{1}{W} \cdot \sum_{i>W} (f_i^*)^2} \leq \frac{1}{W}.$$

406 To this end, we first show that $f_W^* \leq \frac{1}{W}$. If not, we must have for $i = 1, \dots, W$, $f_i^* \geq f_W^* > \frac{1}{W}$,
 407 as $(f_i^*)_{i=1}^d$ is sorted in non-increasing order. Then $\sum_i f_i^* \geq \sum_{i=1}^W f_i^* > 1$, which contradicts to the
 408 fact that $(f_i^*)_{i=1}^d$ is a frequency vector. We have shown that $f_W^* \leq \frac{1}{W}$, and this further implies that
 409 for any $i \geq W$, $f_i^* \leq f_W^* \leq \frac{1}{W}$. Then we can obtain

$$\sqrt{\frac{1}{W} \cdot \sum_{i>W} (f_i^*)^2} \leq \sqrt{\frac{1}{W^2} \cdot \sum_{i>W} f_i^*} \leq \frac{1}{W},$$

410 since $(f_i^*)_{i=1}^d$ is a frequency vector. We have completed all the proof. \square

411 **C.3 Proof of Corollary 2.3**

412 *Proof of Corollary 2.3.* Define

$$E(W) := \sqrt{\frac{1}{W} \sum_{i>W} (f_i^*)^2}.$$

413 We will show the following:

- 414 1. If $W \geq \#\{f_i \geq \tau\} + \frac{1}{\tau^2} \sum_{f_i < \tau} f_i^2$, then $E(W) \leq \tau$.
- 415 2. Moreover, if $E(W) \leq \tau$, then $W \geq \frac{1}{2}(\#\{f_i \geq \tau\} + \frac{1}{\tau^2} \sum_{f_i < \tau} f_i^2)$.

416 Then Corollary 2.3 follows by combining Corollary 2.2 with the above claims.

417 We first show the first part. First note that $W \geq \#\{f_i \geq \tau\}$ and that $(f_i^*)_{i=1}^d$ is sorted in non-
 418 increasing order, so for all $i \geq W$ it holds that $f_i^* < \tau$. Therefore,

$$E(W) := \sqrt{\frac{1}{W} \sum_{i>W} (f_i^*)^2} \leq \sqrt{\frac{1}{W} \sum_{f_i < \tau} f_i^2}.$$

419 Moreover, note that $W \geq \frac{1}{\tau^2} \sum_{f_i < \tau} f_i^2$, so we further have $E(W) \leq \tau$.

420 To show that second part, we first note that, by definition, $E(W) \leq \tau$ is equivalent to

$$2W \geq W + \frac{1}{\tau^2} \sum_{i>W} (f_i^*)^2.$$

421 Consider the following function

$$F(k) := k + \frac{1}{\tau^2} \sum_{i>k} (f_i^*)^2, \quad k \geq 1,$$

422 one can directly verify that $F(k)$ is minimized at $k^* := \#\{i : f_i \geq \tau\}$; moreover,

$$F(k^*) = k^* + \frac{1}{\tau^2} \sum_{i>k^*} (f_i^*)^2 = \#\{f_i \geq \tau\} + \frac{1}{\tau^2} \sum_{f_i < \tau} f_i^2.$$

423 Therefore, we have

$$2W \geq F(W) \geq F(k^*) = \#\{f_i \geq \tau\} + \frac{1}{\tau^2} \sum_{f_i < \tau} f_i^2.$$

424 This completes our proof. □

425 **D Missing Proofs for Section 3**

426 **D.1 Proof of Theorem 3.1**

427 *Proof of Theorem 3.1.* The proof is motivated by Huang et al. [12].

428 Define the following events

$$E_j^{(m)} := \left\{ |\hat{f}_j^{(m)} - f_j^{(m)}| \leq C \cdot \sqrt{\frac{\log(1/p)}{L} \cdot \frac{1}{W} \cdot \sum_{i>W} (f_i^{(m)})^2} \right\}, m \in [M], j \in [d].$$

429 Then by Proposition 2.1 we have

$$\mathbb{P}\{E_j^{(m)}\} \geq 1 - p.$$

430 Then by union bound, we have

$$\mathbb{P}\left\{ \bigcap_{m=1}^M E_j^{(m)} \right\} \geq 1 - Mp.$$

431 Conditional on the event of $\bigcap_{m=1}^M E_j^{(m)}$, we know that every random variable $\hat{f}_j^{(m)} - f_j^{(m)}$ is bounded
432 within

$$(-F^{(m)}, F^{(m)}),$$

433 where

$$F^{(m)} := C \cdot \sqrt{\frac{\log(1/p)}{L} \cdot \frac{1}{W} \cdot \sum_{i>W} (f_i^{(m)})^2}.$$

434 So by Hoeffding inequality, we have

$$\mathbb{P}\left\{ \left| \frac{1}{M} \sum_{m=1}^M \hat{f}_j^{(m)} - \frac{1}{M} \sum_{m=1}^M f_j^{(m)} \right| \leq \sqrt{\frac{\log(2/p_1)}{2M^2} \sum_{m=1}^M (F^{(m)})^2} \mid \bigcap_{m=1}^M E_j^{(m)} \right\} \geq 1 - p_1$$

435 Then we have

$$\mathbb{P}\left\{ \left| \frac{1}{M} \sum_{m=1}^M \hat{f}_j^{(m)} - \frac{1}{M} \sum_{m=1}^M f_j^{(m)} \right| \leq \sqrt{\frac{\log(2/p_1)}{2M^2} \sum_{m=1}^M (F^{(m)})^2} \right\} \geq 1 - p_1 - Mp.$$

436 Note that

$$\begin{aligned} \frac{\log(2/p_1)}{2M^2} \sum_{m=1}^M (F^{(m)})^2 &= \frac{\log(2/p_1)}{2M^2} \sum_{m=1}^M C^2 \cdot \frac{\log(1/p)}{L} \cdot \frac{1}{W} \cdot \sum_{i>W} (f_i^{(m)})^2 \\ &= C^2 \cdot \frac{\log(2/p_1) \log(1/p)}{2L} \cdot \frac{1}{W} \cdot \sum_{i>W} (F_i)^2 \end{aligned}$$

437 So we have

$$\begin{aligned} \mathbb{P}\left\{ \left| \frac{1}{M} \sum_{m=1}^M \hat{f}_j^{(m)} - \frac{1}{M} \sum_{m=1}^M f_j^{(m)} \right| \leq \sqrt{C^2 \cdot \frac{\log(2/p_1) \log(1/p)}{2L} \cdot \frac{1}{W} \cdot \sum_{i>W} (F_i)^2} \right\} \\ \geq 1 - p_1 - Mp. \end{aligned}$$

438 Note replace $p_1 = p'/2$ and $p = p'/(2M)$, we have that

$$\mathbb{P}\left\{ \left| \frac{1}{M} \sum_{m=1}^M \hat{f}_j^{(m)} - \frac{1}{M} \sum_{m=1}^M f_j^{(m)} \right| \leq \sqrt{C^2 \cdot \frac{\log(1/p') \log(M/p')}{L} \cdot \frac{1}{W} \cdot \sum_{i>W} (F_i)^2} \right\} \geq 1 - p'.$$

439

□

440 **D.2 Proof of Theorem 3.2**

441 *Proof of Theorem 3.2.* Let us consider the hybrid sketch approach in Algorithm 1. Recall that within
 442 a round, clients use the same set of hash functions to construct their sketching matrices. Across
 443 different rounds, clients use the same set of location hashes but a fresh set of sign hashes. Denote the
 444 hash functions by:

$$\begin{aligned} h_\ell &: [d] \rightarrow [w], \quad \ell = 1, \dots, L; \\ \sigma_\ell^{(m)} &: [d] \rightarrow \{+1, -1\}, \quad \ell = 1, \dots, L; m = 1, \dots, M. \end{aligned}$$

445 Recall the local frequency in each round is defined by

$$\mathbf{f}^{(m)} := \frac{1}{n} \sum_{t=1}^n \mathbf{x}^{(m,t)}, \quad m = 1, \dots, M.$$

446 And the global frequency vector is defined by

$$\mathbf{f} := \frac{1}{M} \sum_{m=1}^M \mathbf{f}^{(m)}.$$

447 Then according to the communication protocol, the server receives M sketching matrices (each
 448 corresponds to a summation of clients' sketches within the same round). From the m -th sketch, we
 449 can extract L estimators for each index $j \in [d]$, i.e.,

$$\begin{aligned} \tilde{\mathbf{f}}_j^{(m,\ell)} &:= \sum_{i=1}^d \mathbb{1}[h_\ell(i) = h_\ell(j)] \cdot \sigma_\ell^{(m)}(j) \cdot \sigma_\ell^{(m)}(i) \cdot \mathbf{f}_i^{(m)}, \quad j \in [d], m \in [M], \ell \in [L] \\ &= \mathbf{f}_j^{(m)} + \sum_{i \neq j} \mathbb{1}[h_\ell(i) = h_\ell(j)] \cdot \sigma_\ell^{(m)}(j) \cdot \sigma_\ell^{(m)}(i) \cdot \mathbf{f}_i^{(m)}. \end{aligned}$$

450 For each index, we will first average the estimators from different rounds to reduce the variance,
 451 then take the median over different rows to amplify the success probability. In particular, denote the
 452 round-wise averaging by

$$\begin{aligned} \tilde{\mathbf{f}}_j^{(\ell)} &:= \frac{1}{M} \sum_{m=1}^M \tilde{\mathbf{f}}_j^{(m,\ell)}, \quad j \in [d], \ell \in [L] \\ &= \frac{1}{M} \sum_{m=1}^M \mathbf{f}_j^{(m)} + \frac{1}{M} \sum_{m=1}^M \sum_{i \neq j} \mathbb{1}[h_\ell(i) = h_\ell(j)] \cdot \sigma_\ell^{(m)}(j) \cdot \sigma_\ell^{(m)}(i) \cdot \mathbf{f}_i^{(m)} \\ &= \underbrace{\mathbf{f}_j}_{\text{signal}} + \underbrace{\frac{1}{M} \sum_{i \neq j} \mathbb{1}[h_\ell(i) = h_\ell(j)] \cdot \sum_{m=1}^M \sigma_\ell^{(m)}(j) \cdot \sigma_\ell^{(m)}(i) \cdot \mathbf{f}_i^{(m)}}_{\text{noise}} \\ &= \underbrace{\mathbf{f}_j}_{\text{signal}} + \underbrace{\frac{1}{M} \sum_{i \neq j, i \in \mathbb{W}} \mathbb{1}[h_\ell(i) = h_\ell(j)] \cdot \sum_{m=1}^M \sigma_\ell^{(m)}(j) \cdot \sigma_\ell^{(m)}(i) \cdot \mathbf{f}_i^{(m)}}_{\text{headNoise}} \\ &\quad + \underbrace{\frac{1}{M} \sum_{i \neq j, i \notin \mathbb{W}} \mathbb{1}[h_\ell(i) = h_\ell(j)] \cdot \sum_{m=1}^M \sigma_\ell^{(m)}(j) \cdot \sigma_\ell^{(m)}(i) \cdot \mathbf{f}_i^{(m)}}_{\text{tailNoise}}. \end{aligned} \tag{5}$$

453 Then we take the median over these estimators to obtain

$$\tilde{\mathbf{f}}_j := \text{median}\{\tilde{\mathbf{f}}_j^{(\ell)}, \ell \in [L]\}, \quad j \in [d].$$

454 **Head Noise.** The only randomness comes from the algorithm. Note that the head noise contains at
 455 most $|\mathbb{W}| \leq 0.1W$ independent terms, and each is zero with probability $1 - 1/W$. Thus the head
 456 noise is zero with probability at least $(1 - 1/W)^{|\mathbb{W}|} \geq (1 - 1/W)^{0.1W} \geq 0.9$ provided that $W > 10$.

457 **Tail Noise.** Now consider the second noise term in (5). Fixing ℓ and j . Define

$$\begin{aligned}\xi_i^{(m)} &:= \sigma_\ell^{(m)}(j) \cdot \sigma_\ell^{(m)}(i) \cdot \mathbf{f}_i^{(m)} \\ \xi_i &:= \sum_{m=1}^M \xi_i^{(m)} = \sum_{m=1}^M \sigma_\ell^{(m)}(j) \cdot \sigma_\ell^{(m)}(i) \cdot \mathbf{f}_i^{(m)} \\ \eta_i &:= \mathbb{1}[h(i) = h(j)] \\ \text{tailNoise} &:= \frac{1}{M} \sum_{i \neq j, i \notin \mathbb{W}} \eta_i \cdot \xi_i.\end{aligned}$$

458 First notice that $(\xi_i^{(m)})_{m=1}^M$ are independent random variables and

$$\mathbb{E}[\xi_i^{(m)}] = 0, \quad \text{Var}[\xi_i^{(m)}] = (\mathbf{f}_i^{(m)})^2.$$

459 These imply that

$$\mathbb{E}[\xi_i] = 0, \quad \text{Var}[\xi_i] = \sum_{m=1}^M (\mathbf{f}_i^{(m)})^2.$$

460 Moreover, notice that $(\eta_i, \xi_i)_{i \neq j}$ are independent random variables, and

$$\mathbb{E}[\eta_i^2] = \frac{1}{W},$$

461 we then have

$$\begin{aligned}\mathbb{E}[\eta_i \xi_i] &= 0; \\ \text{Var}[\eta_i \xi_i] &= \mathbb{E}[\eta_i^2] \cdot \text{Var}[\xi_i] + \text{Var}[\eta_i] \cdot (\mathbb{E}[\xi_i])^2 \\ &= \frac{1}{W} \cdot \sum_{m=1}^M (\mathbf{f}_i^{(m)})^2.\end{aligned}$$

462 Therefore we conclude that

$$\begin{aligned}\mathbb{E}[\text{tailNoise}] &= \frac{1}{M} \sum_{i \neq j, i \notin \mathbb{W}} \mathbb{E}[\eta_i \xi_i] = 0; \\ \text{Var}[\text{tailNoise}] &= \frac{1}{M^2} \sum_{i \neq j, i \notin \mathbb{W}} \text{Var}[\eta_i \xi_i] \\ &= \frac{1}{M^2 W} \cdot \sum_{i \neq j, i \notin \mathbb{W}} \sum_{m=1}^M (\mathbf{f}_i^{(m)})^2 \\ &\leq \frac{1}{M^2 W} \cdot \sum_{i \notin \mathbb{W}} \sum_{m=1}^M (\mathbf{f}_i^{(m)})^2.\end{aligned}$$

463 Then by Chebyshev we see that: for fixed $j \in [d]$ and $\ell \in [L]$ it holds that

$$\mathbb{P}\left\{|\text{tailNoise}| \geq \sqrt{\frac{10}{M^2 W} \cdot \sum_{i \notin \mathbb{W}} \sum_{m=1}^M (\mathbf{f}_i^{(m)})^2}\right\} < 0.1.$$

464 By a union bound we see that: for fixed $j \in [d]$ and $\ell \in [L]$ it holds that

$$\mathbb{P}\left\{|\tilde{\mathbf{f}}_j^{(\ell)} - \mathbf{f}_j| < \sqrt{\frac{10}{M^2 W} \cdot \sum_{i \notin \mathbb{W}} \sum_{m=1}^M (\mathbf{f}_i^{(m)})^2}\right\} > 0.8 > 0.5.$$

465 **Probability Amplification.** Fixing j . Recall that $(\tilde{\mathbf{f}}_j^{(\ell)})_{\ell=1}^L$ are i.i.d. random variables and that
 466 $\tilde{\mathbf{f}}_j := \text{median}\{\tilde{\mathbf{f}}_j^{(\ell)} : \ell \in [L]\}$. By Chernoff over ℓ and union bound over j we see that:

$$\mathbb{P}\left\{\text{for each } j \in [d], \quad |\tilde{\mathbf{f}}_j - \mathbf{f}_j| \geq \sqrt{\frac{10}{M^2 W} \cdot \sum_{i \notin \mathbb{W}} \sum_{m=1}^M (\mathbf{f}_i^{(m)})^2}\right\} < 2d \cdot \exp(-\Omega(L)).$$

467 By choosing $L = \Theta(\log(2d/\delta))$ we obtain that, with probability at least $1 - \delta$,

$$\text{for each } j \in [d], \quad |\tilde{\mathbf{f}}_j - \mathbf{f}_j| \lesssim \sqrt{\frac{10}{M^2 W} \cdot \sum_{i \notin \mathbb{W}} \sum_{m=1}^M (\mathbf{f}_i^{(m)})^2}.$$

468

□

469 **E Missing Proofs for Section 4**

470 **E.1 Proof of Theorem 4.1**

471 *Proof of Theorem 4.1.* We follow the method of Pagh and Thorup [14], Zhao et al. [17] to add DP
 472 noise to all M sketches. Suppose $\mathcal{F} = (\mathbf{f}^{(m)})_{m=1}^M$ and $\mathring{\mathcal{F}} = (\mathring{\mathbf{f}}^{(m)})_{m=1}^M$ are the sets of local
 473 frequencies for two neighboring datasets respectively, then

$$\|\mathcal{F} - \mathring{\mathcal{F}}\|_2 \leq \frac{1}{n}.$$

474 Denote the sketches to be released by $\mathcal{S} \circ \mathcal{F} := (\mathbf{S}^{(m)} \circ \mathbf{f}^{(m)})_{m=1}^M$. One can then calculate the
 475 ℓ_2 -sensitivity:

$$\|\mathcal{S} \circ \mathcal{F} - \mathcal{S} \circ \mathring{\mathcal{F}}\|_2 \leq \frac{\sqrt{L}}{n},$$

476 where $L \approx \log(d/\delta)$ is the sketch length. Therefore the sketching will be (ϵ, δ) -DP by adding
 477 Gaussian noise $\mathcal{N}(0, \sigma^2)$ to each bucket of each sketch, where

$$\sigma \approx \frac{\sqrt{L \log(1/\delta)}}{n\epsilon}.$$

478 The final released frequency estimator is obtained by post-processing the sketch, so it is also (ϵ, δ) -DP.
 479 We then calculate the error for the noisy sketch matrix. For each row estimator, we have that with
 480 probability at least $2/3$:

$$\begin{aligned} \tilde{\mathbf{f}}_j^{(\ell)} - \mathbf{f}_j^\ell &= \text{tailNoise} + \frac{1}{M} \sum_{m=1}^M \text{rad}_m \cdot \mathcal{N}(0, \sigma^2) \\ &= \text{tailNoise} + \mathcal{N}(0, \sigma^2/M) \\ &\lesssim \sqrt{\frac{1}{M^2 w} \cdot \sum_{i \notin \mathbb{W}} \sum_{m=1}^M (\mathbf{f}_i^{(m)})^2} + \frac{\sqrt{L \log(1/\delta)}}{\sqrt{M} n \epsilon}. \end{aligned}$$

481 By taking median over $L \approx \log(d/\delta)$ repeats, we see that with probability at least $1 - \delta$, it holds that

$$\text{for each } j \in [d], \quad |\hat{\mathbf{f}}_j - \mathbf{f}_j| \lesssim \sqrt{\frac{1}{M^2 w} \cdot \sum_{i \notin \mathbb{W}} \sum_{m=1}^M (\mathbf{f}_i^{(m)})^2} + \frac{\sqrt{\log(d/\delta) \cdot \log(1/\delta)}}{\sqrt{M} n \epsilon}.$$

482

□