

A Locally Differential Private Coding-Assisted Succinct Histogram Protocol

Hsuan-Po Liu* and Hessam MahdaviFar*[†]

*Department of Electrical Engineering and Computer Science, University of Michigan, Ann Arbor, MI 48109, USA

[†]Department of Electrical and Computer Engineering, Northeastern University, Boston, MA 02115, USA

Emails: hsuampo@umich.edu, h.mahdaviFar@northeastern.edu

Abstract—A succinct histogram captures frequent items and their frequencies across clients and has become increasingly important for large-scale, privacy-sensitive machine learning applications. To develop a rigorous framework to guarantee privacy for the succinct histogram problem, local differential privacy (LDP) has been utilized and shown promising results. To preserve data utility under LDP, which essentially works by intentionally adding noise to data, error-correcting codes naturally emerge as a promising tool for reliable information collection. This work presents the first practical (ϵ, δ) -LDP protocol for constructing succinct histograms using error-correcting codes. To this end, polar codes and their successive-cancellation list (SCL) decoding algorithms are leveraged as the underlying coding scheme. More specifically, our protocol introduces Gaussian-based perturbations to enable efficient soft decoding. Experiments demonstrate that our approach outperforms prior methods, particularly for items with low true frequencies, while maintaining similar frequency estimation accuracy.

Index Terms—Succinct histogram, local differential privacy, error-correcting codes.

I. INTRODUCTION

In the era of big data, collecting large-scale contextual information is key to generating accurate statistics, training sophisticated machine learning models, and turning raw data into actionable insights, which has been well-studied for decades [1]. However, the state-of-the-art protocols that directly gather user details raise serious privacy concerns. Data on local devices often includes sensitive personal information, such as behavioral traits, usage patterns, or private preferences, which raises substantial privacy concerns when raw records are collected without protection. By adopting privacy-preserving techniques [2], such as perturbations on the information shared by the devices, we can enable powerful analytics and learning workflows that, otherwise, would be impractical or impossible under traditional data collection paradigms.

Differential Privacy (DP) [2], [3] provides a mathematically rigorous framework for privacy-preserving data analysis. It ensures that the addition or removal of any single individual’s record has a negligible effect on the output of an algorithm, thereby limiting the risk of sensitive information leakage. In the centralized DP model, a trusted curator aggregates raw data and applies this noise before releasing outputs. However, such trustworthiness assumptions may not hold in many practical settings, simply due to the large scale of the systems and the volume of data that are collected from a large number of clients [4]. Local DP (LDP) [5], on the other hand, eliminates the need

for a trusted intermediary by shifting the perturbation step to each user’s device: individuals perturb their data locally before transmission. Therefore, the LDP framework better meets the needs of modern large-scale data collection.

A *succinct histogram* lists heavy hitters and estimates their frequencies [6]. Coding-based methods [6]–[8] encode client data with error-correcting codes, apply random perturbations to satisfy ϵ -LDP, and decode at the server to identify heavy hitters. Early work [6] introduced a protocol, referred to as $\text{PROT}^{\text{PP}}\text{-S-Hist}_{\text{PP}}$, to solve the *unique heavy hitter* problem, aiming to identify items held by at least an η fraction of clients. Although $\text{PROT}^{\text{PP}}\text{-S-Hist}_{\text{PP}}$ provided error bounds, it lacked practical code design considerations and operated on hard decisions. Later refinements in [7], [8] improved the theoretical error bounds, but remained untested in practice.

This paper proposes the first practical (ϵ, δ) -LDP protocol for succinct histograms that utilizes error-correcting codes, in particular, polar codes [9] and their successive-cancellation list (SCL) decoder [10]. By adopting the analytic Gaussian mechanism [11] for perturbation, the protocol enables efficient soft decoding. This is appealing from a practical perspective, indicating that existing mechanisms for reliable communication currently deployed, e.g., 5G protocols, can be re-used to guarantee privacy as well without incurring additional computational cost [12]. We provide theoretical bounds on the frequency estimation error and extend the design to the general succinct histogram setting [6]. For comparison, we implement $\text{PROT}^{\text{PP}}\text{-S-Hist}_{\text{PP}}$ with polar codes and maximum likelihood decoding. Experiments show that our protocol achieves a lower threshold on η , offering greater robustness in decoding errors while maintaining a comparable frequency estimation accuracy. The main contributions are: (i) the first practical (ϵ, δ) -LDP coding-assisted succinct histogram protocol; (ii) leveraging analytic Gaussian mechanism to enable soft decoding; and (iii) experimental validation demonstrating improved robustness over prior work.

The rest of the paper is structured as follows. In Section II, some preliminaries are provided. In Section III, we present the proposed protocol with (ϵ, δ) -LDP guarantees, followed by the analysis of frequency estimation error. The simulation results are included in Section IV. Finally, we conclude the paper in Section V.

II. PRELIMINARIES

A. Succinct histogram

The succinct histogram problems analyze the most frequent items, the heavy hitters, from clients exceeding a certain frequency threshold. We define a succinct histogram and the unique heavy hitter problem we consider in this paper. The details can be found in [6].

Definition 1 (Succinct Histogram). A succinct histogram is a data structure that provides a list of frequent items called heavy hitters, together with their corresponding estimated frequencies among the clients.

Definition 2 (The unique heavy hitter problem). At least η fraction of clients hold the same item as an k -bit binary vector \mathbf{u}^* for some $\mathbf{u}^* \in \mathcal{U} = \{0, 1\}^k$ unknown to the untrusted server, while other clients hold a special symbol \perp , to represent "no item." The goal is to design an efficient (ϵ, δ) -LDP protocol of a succinct histogram, for η as small as possible.

B. Local Differential Privacy

The main idea to achieve (ϵ, δ) -LDP is to intentionally add perturbation via random noise sampled from a probability distribution at the client's side, before sharing information with the untrusted server. The definition of (ϵ, δ) -LDP is provided below. The details can be found in [3], [5].

Definition 3 $((\epsilon, \delta)$ -LDP). Let u and u' be neighboring datasets that differ only by a single event, i.e., $\text{dist}(u, u') = 1$. The neighboring datasets u and u' satisfy (ϵ, δ) -LDP for any $\epsilon > 0$ under a randomized mechanism \mathcal{R} that under $\mathcal{E} \subseteq \text{Range}(\mathcal{R})$,

$$\mathbb{P}[\mathcal{R}(u) \in \mathcal{E}] \leq e^\epsilon \cdot \mathbb{P}[\mathcal{R}(u') \in \mathcal{E}] + \delta, \quad (1)$$

where δ represents the failure probability.

The sensitivity for a query function $q(\cdot)$ is the maximum difference between two queries.

Definition 4 (Sensitivity). For two neighboring datasets u and u' in an individual client together with a query function $q : \mathcal{D} \rightarrow \mathbb{R}$, the sensitivity is defined as follows:

$$\Delta \stackrel{\text{def}}{=} \max_{\text{dist}(u, u')=1} \|q(u) - q(u')\|_2. \quad (2)$$

The Gaussian mechanism is also defined as follows.

Definition 5 (Gaussian mechanism). Consider applying a query function q on a dataset u . Then, the Gaussian mechanism \mathcal{R} is defined as

$$\mathcal{R}(u) \stackrel{\text{def}}{=} q(u) + \mathcal{N}(0, \sigma^2), \quad (3)$$

which adds random noise to the query result according to a zero-mean Gaussian distribution with variance σ^2 . Note that σ should be a function of ϵ , δ , and Δ .

Intuitively, σ should be proportional to sensitivity Δ . That is, a greater value of sensitivity for a dataset implies that we need a larger perturbation of the query, making it more indistinguishable from the adversary. We indicate the expression

of σ for the randomized mechanism in our proposed protocol in the next section.

C. Polar Codes

a) *Coding structure*: An (n, k) -Polar code [9] is a linear block code that has a code dimension of k and the code length as $n = 2^m$, where m is the Kronecker power such that the generator matrix \mathbf{G} follows the transformation $\mathbf{G} = \mathbf{G}_2^{\otimes m}$ given the base matrix $\mathbf{G}_2 = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$. We denote the binary information vector by \mathbf{u} and the transmitted vector as a codeword corresponding to each information vector by \mathbf{x} , such that $\mathbf{x} = \mathbf{u}\mathbf{G}$, where $\mathbf{x} \in \{0, 1\}^n$ and $\mathbf{u} \in \{0, 1\}^k$. The key idea behind polar codes is the phenomenon of channel polarization, where certain bit-channels become highly reliable while others become completely noisy as the code length increases. A polar encoder exploits this polarization effect to transmit information bits through reliable bit-channels, while fixing the values of the less reliable ones to zeros, called frozen bits.

b) *Decoder*: The SC decoder [9] for polar codes operates by sequentially estimating each bit in the codeword using a recursive computation of log-likelihood ratios (LLRs). To estimate the j th message bit u_j in \mathbf{u} , the SC decoder makes the following decision based on the received vector \mathbf{y} and past decisions $\hat{\mathbf{u}}_1^{j-1} = [\hat{u}_1, \dots, \hat{u}_{j-1}]$ as

$$\hat{u}_i = \begin{cases} 0 & \text{if } L^{(j)}(\mathbf{y}, \hat{\mathbf{u}}_1^{j-1}) \geq 1 \\ 1 & \text{otherwise} \end{cases} \quad (4)$$

where $L^{(j)}(\mathbf{y}, \hat{\mathbf{u}}_1^{j-1}) \stackrel{\text{def}}{=} \frac{W^{(j)}(\mathbf{y}, \hat{\mathbf{u}}_1^{j-1}|0)}{W^{(j)}(\mathbf{y}, \hat{\mathbf{u}}_1^{j-1}|1)}$ given that $W^{(j)}(\mathbf{y}, \hat{\mathbf{u}}_1^{j-1}|u_j) = \sum_{\mathbf{u}_{j+1}^{n-1} \in \{0, 1\}^{n-i}} W(\mathbf{y}|\mathbf{x})$. SC decoding remains attractive due to its low complexity and suitability for hardware implementations. It also serves as the foundation for more powerful variants such as the SCL decoder, which is our focus in the proposed protocol.

III. THE PROPOSED (ϵ, δ) -LDP CODING-ASSISTED SUCCINCT HISTOGRAM PROTOCOL

In this section, we first illustrate our proposed protocol. Next, we analyze the (ϵ, δ) -LDP guarantees, followed by the frequency estimation error analysis.

A. The Proposed Protocol

We consider N clients with a single untrusted server that collects information from all clients for the unique heavy hitter problem in Definition 2. Client i either holds the true item as a k -bit binary vector $\mathbf{u}^* \in \mathcal{U} = \{0, 1\}^k$ or has no item as \perp , denoted by $\mathbf{u}_i \in \{\mathbf{u}^*\} \cup \{\perp\}$, for $i \in \{1, \dots, N\} \stackrel{\text{def}}{=} [N]$. We denote the set of clients that hold the unique item \mathbf{u}^* as \mathcal{T} , so we have

$$\mathbf{u}_i = \begin{cases} \mathbf{u}^* & \text{if } i \in \mathcal{T} \\ \perp & \text{if } i \in [N] \setminus \mathcal{T} \end{cases} \quad (5)$$

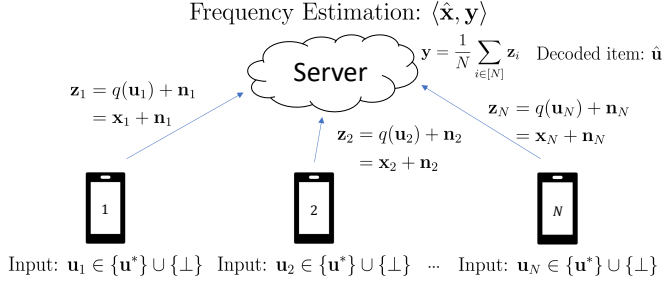


Fig. 1: Proposed Coding-Assisted Succinct Histogram Protocol

for $i \in [N]$. According to Definition 2, at least η fraction of the clients have the unique item \mathbf{u}^* . Thus, we define the *true frequency* of such item as $f(\mathbf{u}^*) \stackrel{\text{def}}{=} \frac{|\mathcal{T}|}{N} \geq \eta$.

Clients $i \in \mathcal{T}$ who hold the unique item $\mathbf{u}_i = \mathbf{u}^*$ encode \mathbf{u}_i by the given (n, k) -polar code with a generator matrix \mathbf{G} to obtain the corresponding codeword $\mathbf{v}_i = \mathbf{u}_i \mathbf{G} = \mathbf{u}^* \mathbf{G}$, then modulate the codeword by a normalized BPSK, which has the mapping $\{0, 1\}^n \rightarrow \{-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}\}^n$, to generate $q(\mathbf{u}_i) = \mathbf{x}_i = \frac{1}{\sqrt{n}}(2\mathbf{v}_i - \mathbf{1}_n) = \frac{1}{\sqrt{n}}(2\mathbf{u}^* \mathbf{G} - \mathbf{1}_n)$; for clients $i \in [N] \setminus \mathcal{T}$ who do not have the unique item, i.e. $\mathbf{u}_i = \perp$, set $q(\mathbf{u}_i) = \mathbf{x}_i = \mathbf{0}_n$. Thus, we summarize $q(\mathbf{u}_i) = \mathbf{x}_i$ as follows:

$$q(\mathbf{u}_i) = \mathbf{x}_i = \begin{cases} \frac{1}{\sqrt{n}}(2\mathbf{u}^* \mathbf{G} - \mathbf{1}_n) & \text{if } \mathbf{u}_i = \mathbf{u}^* \\ \mathbf{0}_n & \text{if } \mathbf{u}_i = \perp \end{cases} \quad (6)$$

for $i \in [N]$. Next, client i applies the randomized mechanism \mathcal{R}_i , that satisfies (ϵ, δ) -LDP, to obtain $\mathbf{z}_i = \mathcal{R}_i(\mathbf{u}_i) \stackrel{\text{def}}{=} q(\mathbf{u}_i) + \mathbf{n}_i = \mathbf{x}_i + \mathbf{n}_i$, where \mathbf{n}_i is a Gaussian noise vector that all entries are sampled from a zero mean Gaussian distribution $\mathcal{N}(0, \sigma^2)$, then sends \mathbf{z}_i to the server, for $i \in [N]$. We leave the detailed setup of the randomized mechanism \mathcal{R}_i to the next subsection.

The server collects \mathbf{z}_i 's from clients $i \in [N]$, then computes the mean $\mathbf{y} = \frac{1}{N} \sum_{i \in [N]} \mathbf{z}_i$. An SCL decoder is applied to decode the estimated item $\hat{\mathbf{u}}$ based on the LLR $\hat{\mathbf{y}} = \frac{2}{\sigma^2} \mathbf{y}$ such that $\hat{\mathbf{u}} = \text{Dec}(\hat{\mathbf{y}})$. The frequency is estimated by calculating $\hat{f}(\hat{\mathbf{u}}) = \langle q(\hat{\mathbf{u}}), \mathbf{y} \rangle = \langle \hat{\mathbf{x}}, \mathbf{y} \rangle$, where $q(\hat{\mathbf{u}}) = \hat{\mathbf{x}} = \frac{1}{\sqrt{n}}(2\hat{\mathbf{v}} - \mathbf{1}_n) = \frac{1}{\sqrt{n}}(2\hat{\mathbf{u}} \mathbf{G} - \mathbf{1}_n)$. Finally, the succinct histogram is generated as $(\hat{\mathbf{u}}, \hat{f}(\hat{\mathbf{u}}))$, which includes the estimated item and the estimated frequency, respectively. We summarize the protocol in Fig. 1.

B. The (ϵ, δ) -LDP guarantee

To ensure (ϵ, δ) -LDP by adding perturbations sampled from a zero mean Gaussian distribution $\mathcal{N}(0, \sigma^2)$, we must identify the variance σ^2 of such distributions. We start by analyzing the sensitivity of the proposed protocol as follows.

Lemma 1 (Sensitivity). *The sensitivity of our proposed protocol is bounded as follows:*

$$\Delta = \max_{\mathbf{u}^*, \hat{\mathbf{u}}} \|q(\mathbf{u}^*) - q(\hat{\mathbf{u}})\|_2 = 2. \quad (7)$$

Proof. Given that $q(\mathbf{u}^*) = \mathbf{x}^* = \frac{1}{\sqrt{n}}(2\mathbf{u}^* \mathbf{G} - \mathbf{1}_n)$ and $q(\hat{\mathbf{u}}) = \hat{\mathbf{x}} = \frac{1}{\sqrt{n}}(2\hat{\mathbf{u}} \mathbf{G} - \mathbf{1}_n)$, we have

$$\begin{aligned} \|q(\mathbf{u}^*) - q(\hat{\mathbf{u}})\|_2 &= \left\| \frac{1}{\sqrt{n}}(2\mathbf{u}^* \mathbf{G} - \mathbf{1}_n) - \frac{1}{\sqrt{n}}(2\hat{\mathbf{u}} \mathbf{G} - \mathbf{1}_n) \right\|_2 \\ &= \left\| \frac{2}{\sqrt{n}}(\mathbf{u}^* - \hat{\mathbf{u}}) \mathbf{G} \right\|_2 \\ &\leq \left\| \frac{2}{\sqrt{n}}(\mathbf{1}_k - \mathbf{0}_k) \mathbf{G} \right\|_2 \\ &= \left\| \frac{2}{\sqrt{n}} \cdot \mathbf{1}_n \right\|_2 = 2. \end{aligned} \quad (8)$$

Thus, the maximum value of $\|q(\mathbf{u}^*) - q(\hat{\mathbf{u}})\|_2$ given any $\mathbf{u}^*, \hat{\mathbf{u}} \in \mathcal{U}$, which is the sensitivity, is $\Delta = 2$. Also, this maximum occurs since the all-one vector is a codeword in the polar codebook. \square

Remark 1. We can reduce the sensitivity of our protocol, as stated in (8), by carefully adjusting the code structure. This can be done by simply setting some of the information bits in the given (n, k) -polar code to zeros. For instance, one can exclude the all-one codeword by setting the information bit corresponding to the last index to zero. Thus, the sensitivity becomes $\Delta = \left\| \frac{2}{\sqrt{n}} \cdot \mathbf{c}_{\text{max-weight}} \right\|_2 < 2$, where $\mathbf{c}_{\text{max-weight}}$ is the codeword with the maximum weight excluding the all-one codeword $\mathbf{1}_n$. This gives rise to an interesting problem, i.e., minimizing the maximum weight of codewords, for polar codes which we leave for future work.

With the sensitivity Δ , to ensure (ϵ, δ) -LDP of our protocol, we employ the Analytic Gaussian mechanism [11] that achieves the optimal noise variance. This leads to the following theorem with a proof that follows from the (ϵ, δ) -LDP guarantee of the underlying analytic Gaussian mechanism.

Theorem 2 ((ϵ, δ) -LDP for the proposed protocol). *The protocol is (ϵ, δ) -LDP with the analytic Gaussian mechanism as the randomized mechanism \mathcal{R}_i , such that $\mathcal{R}_i(\mathbf{u}_i) = q(\mathbf{u}_i) + \mathbf{n}_i = \mathbf{x}_i + \mathbf{n}_i$, where \mathbf{n}_i is a noise vector sampled from $\mathcal{N}(0, \sigma^2)$, which is a zero-mean Gaussian distribution with variance $\sigma^2 = \text{ANALYTICGAUSSIAN}(\epsilon, \delta, \Delta)$, for $i \in [N]$.*

With Theorem 2, we summarize our protocol in Algorithm 1.

Remark 2. The major differences between our proposed protocol Algorithm 1 compared to the most relevant prior work $\text{PROT}^{\text{PP}}\text{-S-Hist}_{\text{PP}}$ in [6, Algorithm 2] are as follows: (1) We formulate the practical code construction for encoding the unique item to the codewords in step 2, which $\text{PROT}^{\text{PP}}\text{-S-Hist}_{\text{PP}}$ did not specify the construction at all; (2) The randomizers between our protocol and $\text{PROT}^{\text{PP}}\text{-S-Hist}_{\text{PP}}$ are fundamentally different. $\text{PROT}^{\text{PP}}\text{-S-Hist}_{\text{PP}}$ randomly picks the j th entry from \mathbf{x}_i , which is denoted by $x_{i,j}$, randomizes it to either $c_\epsilon n x_{i,j}$ or $-c_\epsilon n x_{i,j}$ according to the randomized response [13] when $\mathbf{x}_i \neq \mathbf{0}_n$, or assigning $c_\epsilon \sqrt{n}$ or $-c_\epsilon \sqrt{n}$ uniformly at random, while other entries are set to zeros. Our proposed protocol randomizes each entry by adding

Algorithm 1 (ϵ, δ) -LDP Coding-Assisted Succinct Histogram Protocol for the Unique Heavy Hitter Problem

Input: Clients' inputs $\{\mathbf{u}_i \in \mathbf{u}^* \cup \perp : i \in [N]\}$, privacy guarantee ϵ and δ , sensitivity Δ , and generator matrix \mathbf{G} constructed by an (n, k) -polar code

Output: Succinct histogram $(\hat{\mathbf{u}}, \hat{f}(\hat{\mathbf{u}}))$

- 1: **for** Client $i \in [N]$ **do**
 - 2: If $\mathbf{u}_i \neq \perp$, client i encodes then modulates its item $\mathbf{x}_i = \frac{1}{\sqrt{n}}(2\mathbf{v}_i - \mathbf{1}_n) = \frac{1}{\sqrt{n}}(2\mathbf{u}_i\mathbf{G} - \mathbf{1}_n)$. Else, $\mathbf{x}_i = \mathbf{0}_n$.
 - 3: Client i randomized its private report $\mathbf{z}_i = \mathcal{R}_i(\mathbf{x}_i) = \mathbf{x}_i + \mathbf{n}_i$, where \mathbf{n}_i is a vector that all entries are sampled from a zero mean Gaussian distribution $\mathcal{N}(0, \sigma^2)$ with $\sigma^2 = \text{AnalyticGaussian}(\epsilon, \delta, \Delta)$.
 - 4: Client i sends \mathbf{z}_i to the server.
 - 5: **end for**
 - 6: Server collects \mathbf{z}_i 's from clients $i \in [N]$.
 - 7: Server computes $\mathbf{y} = \frac{1}{N} \sum_{i \in [N]} \mathbf{z}_i$.
 - 8: Server calculates the LLR $\tilde{\mathbf{y}} = \frac{2}{\sigma^2} \mathbf{y}$.
 - 9: Server decodes $\tilde{\mathbf{y}}$ into an estimate for the unique item $\hat{\mathbf{u}} = \text{Dec}(\tilde{\mathbf{y}})$ by an SCL decoder.
 - 10: Server encodes then modulates its decoded item $\hat{\mathbf{x}} = \frac{1}{\sqrt{n}}(2\hat{\mathbf{v}} - \mathbf{1}_n) = \frac{1}{\sqrt{n}}(2\hat{\mathbf{u}}\mathbf{G} - \mathbf{1}_n)$.
 - 11: Server computes a frequency estimate $\hat{f}(\hat{\mathbf{u}}) = \langle \hat{\mathbf{x}}, \mathbf{y} \rangle$.
 - 12: **Return** the succinct histogram $(\hat{\mathbf{u}}, \hat{f}(\hat{\mathbf{u}}))$
-

Gaussian noise. (3) Our protocol enables flexibility in the soft-decision process by adding perturbations to the queries by the analytic Gaussian mechanism in step 3, that satisfy (ϵ, δ) -LDP with optimal noise. In $\text{PROT}^{\text{PP}}\text{-S-Hist}_{\text{PP}}$, a local randomizer inspired by randomized response was applied to perturb the queries, making it fundamentally binarized to a hard-decision process; (4) The SCL decoder used in step 9 is a soft decoder, while $\text{PROT}^{\text{PP}}\text{-S-Hist}_{\text{PP}}$ rounds \mathbf{y} in our step 7 to $(-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}})^n$, which is limited to a hard decoder. Section IV demonstrates the significance of the decoder, especially when it comes to a low true frequency.

C. Error Analysis

We analyze the error in the frequency estimation of the correctly decoded item, i.e., $\hat{\mathbf{u}} = \mathbf{u}^*$. The following lemma derives the expression of \mathbf{y} , which is the mean of collected \mathbf{z}_i 's from clients $i \in [N]$.

Lemma 3. *The received vector at the server in step 7 in Algorithm 1 is*

$$\mathbf{y} = \frac{1}{N} \sum_{i \in [N]} \mathbf{z}_i = f(\mathbf{u}^*)\mathbf{x}^* + \frac{1}{N} \sum_{i \in [N]} \mathbf{n}_i, \quad (9)$$

Proof. Given that $\mathbf{z}_i = \mathbf{x}_i + \mathbf{n}_i$ for clients $i \in [N]$, we have

$$\mathbf{y} = \frac{1}{N} \sum_{i \in [N]} \mathbf{z}_i = \frac{1}{N} \sum_{i \in [N]} \mathbf{x}_i + \frac{1}{N} \sum_{i \in [N]} \mathbf{n}_i. \quad (10)$$

Then, note that $\mathbf{x}_i = \mathbf{x}^*$ for $i \in \mathcal{T}$ and $\mathbf{x}_i = \mathbf{0}_i$ for $i \in [N] \setminus \mathcal{T}$.

$$\begin{aligned} \frac{1}{N} \sum_{i \in [N]} \mathbf{x}_i &= \frac{1}{N} \left[\sum_{i \in \mathcal{T}} \mathbf{x}_i + \sum_{i \in [N] \setminus \mathcal{T}} \mathbf{x}_i \right] \\ &= \frac{1}{N} [|\mathcal{T}| \cdot \mathbf{x}^* + (N - |\mathcal{T}|) \cdot \mathbf{0}_n] \\ &= \frac{|\mathcal{T}|}{N} \mathbf{x}^* = f(\mathbf{u}^*)\mathbf{x}^*. \end{aligned} \quad (11)$$

Combining (10) and (11) proves the statement. \square

Lemma 4 (Error analysis for the correctly decoded items). *Assuming the unique item is correctly decoded, i.e., $\hat{\mathbf{u}} = \mathbf{u}^*$, the frequency estimation error is expressed as follows:*

$$\text{Err}_{\hat{\mathbf{u}}, \mathbf{u}^*} \stackrel{\text{def}}{=} \left| \hat{f}(\hat{\mathbf{u}}) - f(\mathbf{u}^*) \right| = \left| \langle \mathbf{x}^*, \frac{1}{N} \sum_{i \in [N]} \mathbf{n}_i \rangle \right|. \quad (12)$$

Proof. When we have the correct decoded item, we can substitute $\hat{\mathbf{u}}$ by \mathbf{u}^* , thus the problem becomes $\text{Err}_{\hat{\mathbf{u}}=\mathbf{u}^*} = \left| \hat{f}(\mathbf{u}^*) - f(\mathbf{u}^*) \right|$. Given the frequency estimate $\hat{f}(\hat{\mathbf{u}}) = \langle \hat{\mathbf{x}}, \mathbf{y} \rangle$ in the protocol, plugging $\hat{\mathbf{u}} = \mathbf{u}^*$, we obtain that $\hat{f}(\mathbf{u}^*) = \langle \mathbf{x}^*, \mathbf{y} \rangle$, where $\hat{\mathbf{x}} = \mathbf{x}^*$, yields $\text{Err}_{\hat{\mathbf{u}}, \mathbf{u}^*} = |\langle \mathbf{x}^*, \mathbf{y} \rangle - f(\mathbf{u}^*)|$. Then, with Lemma 3 we have

$$\begin{aligned} \text{Err}_{\hat{\mathbf{u}}=\mathbf{u}^*} &= \left| \langle \mathbf{x}^*, f(\mathbf{u}^*)\mathbf{x}^* + \frac{1}{N} \sum_{i \in [N]} \mathbf{n}_i \rangle - f(\mathbf{u}^*) \right| \\ &= \left| f(\mathbf{u}^*)\langle \mathbf{x}^*, \mathbf{x}^* \rangle + \langle \mathbf{x}^*, \frac{1}{N} \sum_{i \in [N]} \mathbf{n}_i \rangle - f(\mathbf{u}^*) \right| \\ &= \left| \langle \mathbf{x}^*, \frac{1}{N} \sum_{i \in [N]} \mathbf{n}_i \rangle \right|, \end{aligned} \quad (13)$$

where $\langle \mathbf{x}^*, \mathbf{x}^* \rangle = 1$. \square

IV. SIMULATION RESULTS

In this section, we analyze how true frequencies, codelength, and number of clients, affect the performance. We consider a $(64, 8)$ -polar code, which implies there are 2^8 possible unique items, with the number of clients $N = 1000$, and we analyze the true frequencies $f(\mathbf{u}^*) \in \{0.5, 0.6, 0.7\}$. The simulations vary ϵ and set $\delta = 10^{-4}$. To decode the estimated unique item, we consider the SC list (SCL) decoder [10] with a list size $L = 8$ for our protocol and the maximum likelihood decoder for $\text{PROT}^{\text{PP}}\text{-S-Hist}_{\text{PP}}$ in [6]. The BLER and frequency estimation error comparisons to ϵ are demonstrated in Figs. 2 and 3, respectively.

Figure 2 demonstrates the robustness of the protocols that experience different true frequencies. When the true frequency is $f(\mathbf{u}^*) = 0.7$, both protocols have good results, but $\text{PROT}^{\text{PP}}\text{-S-Hist}_{\text{PP}}$ achieves a lower BLER due to the high performance of the maximum likelihood decoder. When the true frequency becomes $f(\mathbf{u}^*) = 0.6$, our protocol still perform well, while $\text{PROT}^{\text{PP}}\text{-S-Hist}_{\text{PP}}$ starts to suffer an error floor since the protocol rounds the entries in \mathbf{y} to $\frac{1}{\sqrt{n}}$, if $y_j \geq 0$ and

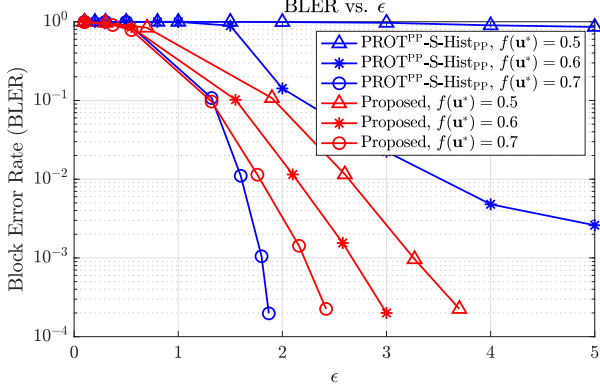


Fig. 2: BLER vs. ϵ for $(n, k) = (64, 8)$ and $N = 1000$

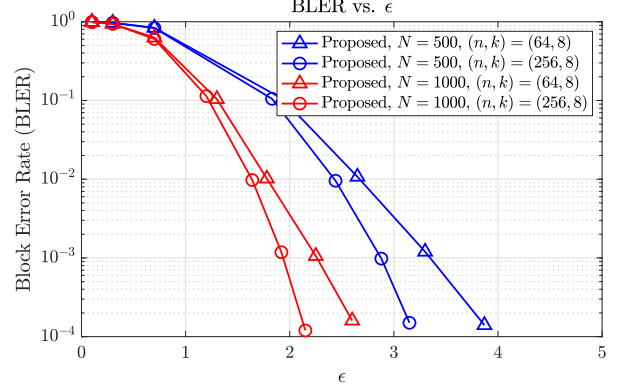


Fig. 4: BLER vs. ϵ for $f(u^*) = 0.7$

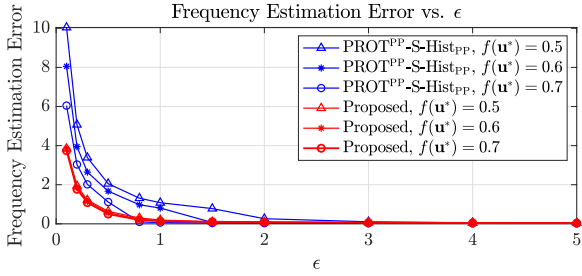


Fig. 3: Frequency Estimation Error vs. ϵ for $(n, k) = (64, 8)$ and $N = 1000$

to $-\frac{1}{\sqrt{n}}$, if $y_j < 0$, for $j \in [n]$. Intuitively, since lower true frequency $f(u^*)$ means more clients have $\mathbf{x}_i = \mathbf{0}_n$, more entries tend to be rounded to $\frac{1}{\sqrt{n}}$, which gradually causes the entries in $\tilde{\mathbf{y}}$ biased towards positive. This becomes worse when the true frequency is decreased to $f(u^*) = 0.5$, where most items are decoded incorrectly. In Fig. 3, one can observe that our frequency estimation outperforms $\text{PROT}^{\text{PP}}\text{-S-Hist}_{\text{PP}}$. These results indicate that our proposed protocol has versatile performance regarding the BLER given a wide range of true frequencies, equivalent to having a high correct decoding rate. Thus, a lower η can be guaranteed such that our protocol is better than $\text{PROT}^{\text{PP}}\text{-S-Hist}_{\text{PP}}$ according to Definition 2, which requires η as small as possible. Furthermore, the proposed protocol has a lower complexity of the decoder while achieving better frequency estimation errors.

In Fig. 4, we evaluate the performance of our proposed protocol in terms of BLER for varying the codelength from $n = 64$ to $n = 256$ and the number of clients from $N = 500$ to $N = 1000$, given $\delta = 10^{-4}$, $f(u^*) = 0.7$, and $k = 8$. We can observe that increasing the codelength from $n = 64$ to $n = 256$ gives our protocol a performance gain for both $N = 500$ and $N = 1000$. Furthermore, increasing the number of clients from $N = 500$ to $N = 1000$ gives us a performance gain, which can be observed from Lemma 4.

V. CONCLUSION

In this work, we introduced the first practically implementable protocol for constructing succinct histograms under

(ϵ, δ) -LDP using error-correcting codes. By leveraging polar codes and an SCL decoder, our design supports soft decoding through an analytic Gaussian mechanism, enabling improved robustness without sacrificing frequency estimation accuracy. Unlike prior works that remained theoretical, we provide concrete constructions and experimental evaluations, demonstrating that our protocol outperforms across varied settings, while achieving similar frequency estimation errors, especially on lower true frequencies of the unique items.

REFERENCES

- [1] J. Leskovec, A. Rajaraman, and J. D. Ullman, *Mining of massive data sets*. Cambridge university press, 2020.
- [2] C. Dwork, F. McSherry, K. Nissim, and A. Smith, “Calibrating noise to sensitivity in private data analysis,” in *Theory of Cryptography: Third Theory of Cryptography Conference, TCC 2006, New York, NY, USA, March 4-7, 2006. Proceedings 3*. Springer, 2006, pp. 265–284.
- [3] C. Dwork, A. Roth et al., “The algorithmic foundations of differential privacy,” *Foundations and Trends® in Theoretical Computer Science*, vol. 9, no. 3–4, pp. 211–407, 2014.
- [4] H.-P. Liu, M. Soleymani, and H. MahdaviFar, “Differentially private coded computing,” in *2023 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2023, pp. 2189–2194.
- [5] J. C. Duchi, M. I. Jordan, and M. J. Wainwright, “Local privacy and statistical minimax rates,” in *2013 IEEE 54th annual symposium on foundations of computer science*. IEEE, 2013, pp. 429–438.
- [6] R. Bassily and A. Smith, “Local, private, efficient protocols for succinct histograms,” in *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*, 2015, pp. 127–135.
- [7] R. Bassily, K. Nissim, U. Stemmer, and A. Thakurta, “Practical locally private heavy hitters,” *Journal of Machine Learning Research*, vol. 21, no. 16, pp. 1–42, 2020.
- [8] M. Bun, J. Nelson, and U. Stemmer, “Heavy hitters and the structure of local privacy,” *ACM Transactions on Algorithms (TALG)*, vol. 15, no. 4, pp. 1–40, 2019.
- [9] E. Arikan, “Channel polarization: A method for constructing capacity-achieving codes for symmetric binary-input memoryless channels,” *IEEE Transactions on information Theory*, vol. 55, no. 7, pp. 3051–3073, 2009.
- [10] I. Tal and A. Vardy, “List decoding of polar codes,” *IEEE transactions on information theory*, vol. 61, no. 5, pp. 2213–2226, 2015.
- [11] B. Balle and Y.-X. Wang, “Improving the gaussian mechanism for differential privacy: Analytical calibration and optimal denoising,” in *International Conference on Machine Learning*, 2018, pp. 394–403.
- [12] H.-P. Liu and H. MahdaviFar, “Projective systematic authentication via reed-muller codes,” in *2024 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2024, pp. 1812–1817.
- [13] S. L. W. and, “Randomized response: A survey technique for eliminating evasive answer bias,” *Journal of the American Statistical Association*, vol. 60, no. 309, pp. 63–69, 1965.