# Communication-Efficient Publication of Sparse Vectors under Differential Privacy

Quentin Hillebrand
The University of Tokyo
Tokyo, Japan
quentin-hillebrand@g.ecc.u-tokyo.ac.jp

Vorapong Suppakitpaisarn
The University of Tokyo
Tokyo, Japan
vorapong@is.s.u-tokyo.ac.jp

Tetsuo Shibuya
The University of Tokyo
Tokyo, Japan
tshibuya@hgc.jp

## ABSTRACT

In this work, we propose a differentially private algorithm for publishing matrices aggregated from sparse vectors. These matrices include social network adjacency matrices, user-item interaction matrices in recommendation systems, and single nucleotide polymorphisms (SNPs) in DNA data. Traditionally, differential privacy in vector collection relies on randomized response, but this approach incurs high communication costs. Specifically, for a matrix with $N$ users, $n$ columns, and $m$ nonzero elements, conventional methods require $\Omega(n \times N)$ communication, making them impractical for large-scale data. Our algorithm significantly reduces this cost to $O(\varepsilon m)$, where $\varepsilon$ is the privacy budget. Notably, this is even lower than the non-private case, which requires $\Omega(m \log n)$ communication. Moreover, as the privacy budget decreases, communication cost further reduces, enabling better privacy with improved efficiency. We theoretically prove that our method yields results identical to those of randomized response, and experimental evaluations confirm its effectiveness in terms of accuracy, communication efficiency, and computational complexity.

## CCS CONCEPTS

• **Security and privacy** → **Privacy-preserving protocols**; **Data anonymization and sanitization**.

## KEYWORDS

Differential privacy, Metric differential privacy, Communication constraint, Graph differential privacy

## 1 INTRODUCTION

**Differential privacy** [4] has emerged as a widely accepted standard for protecting sensitive information while enabling data analysis and sharing. By adding carefully calibrated noise to query results, it ensures that the release of analysis outcomes does not significantly alter an observer's knowledge of users' sensitive information.

To handle cases where data is dispersed across multiple users, various techniques have been developed to generate an obfuscated representation of the raw data [5]. One of the most well-known techniques is **randomized response**. In this technique, each data point is probabilistically altered before being transmitted to the central server. The server then processes the received obfuscated data rather than the original values. This technique is widely recognized for its ability to satisfy local differential privacy [18], a variant of differential privacy that not only safeguards users' information

when publishing analytical results but also ensures privacy during data transmission and storage on the central server.

When a user possesses multiple sensitive data points, a variant of differential privacy known as **metric differential privacy** [1] is commonly employed to enhance privacy protection. This privacy framework enables more accurate data analysis while maintaining user confidentiality. The randomized response technique, which perturbs each data point with a certain probability, also adheres to the requirements of this privacy notion.

We consider the scenario where each user's data is represented as a high-dimensional yet sparse vector—that is, while a user may have a large amount of information, most of the values are identical. For instance, in the context of **graph and social network data**, each user maintains a list of friends, which can be represented as an adjacency vector of size $n$, where $n$ is the total number of users in the system. If a user has only $d \ll n$ friends, the adjacency vector contains just $d$ ones, while the remaining entries, which are zeros (the trivial values in this case), make up nearly the entire vector. A similar situation arises in **recommendation systems**, where each user provides ratings for movies. Since users typically rate only a small subset of all available movies, the majority of entries in their rating vector remain as "N/A," indicating missing or unrated values. In genomic data, it is well established that the vast majority of **genetic information** is identical across all humans, with differences between individuals accounting for only about 1% of the total genetic sequence.

In a non-private setting, collecting and storing sparse vectors is efficient because we only need to track the positions and values of non-trivial elements. Given a vector of size $n$ with at most $d$ non-trivial entries, the communication and storage cost is $O(d \log n)$. However, when applying the randomized response mechanism, each entry is obfuscated with a certain probability, increasing the number of non-trivial elements to $\Omega(n)$ after obfuscation. As a result, the communication and storage cost per user grows to $\Omega(n)$ bits. For a system with $N$ users, the overall communication and storage complexity becomes $\Omega(n \times N)$ bits, which is impractical for many applications.

For example, in a social network setting where $n = N$, this results in a communication and storage cost of $\Omega(N^2)$ bits, making it impractical to process networks with more than approximately $N = 50,000$ nodes. Moreover, some algorithms [16, 11] require the central server to distribute the collected results back to all users, further increasing the overall communication cost to $\Omega(N^3)$, which is prohibitively expensive.

Reducing communication costs has been a focus of many recent studies, including approaches such as rejection sampling [6] and

importance sampling [30]. However, the outputs of these techniques do not precisely match the distribution of the original mechanism. This discrepancy can be problematic, as it may cause the mechanism to lose essential properties such as unbiasedness, making it more challenging to analyze and post-process.

Several mechanisms provide unbiased results [17, 10, 24]. However, the methods proposed in [17, 10] are specifically designed for social network data and still require a communication cost of $\Omega(N^2)$. While the mechanism in [24] is universally applicable, its use for high-dimensional vectors incurs a communication cost of $\Omega(n \times N)$ and suffers from significant computational overhead. In particular, the computational cost grows exponentially with $d$, making it impractical even when the number of non-trivial entries is below 10,000, as the computation time becomes prohibitively high.

### 1.1 Our Contributions

In this paper, we propose a communication- and computation-efficient mechanism for publishing the randomized response of high-dimensional sparse vectors.

Our mechanism incorporates some ideas from the approach called Poisson Private Representation (PPR) in [24]. However, to mitigate the prohibitive computational cost that arises when the number of non-trivial elements $d$ is large, we introduce a random partitioning strategy. Specifically, we divide the high-dimensional vector into $\Theta(d)$ chunks, ensuring that the expected number of non-trivial elements per chunk is $\Theta(1)$. This approach minimizes the likelihood of any chunk containing a disproportionately large number of non-trivial elements, which causes high computation cost in [24]. We then propose an algorithm to efficiently compress the randomized response results for vectors with a small number of non-trivial elements. The complete mechanism is presented in Section 3.

In [24], the algorithm known as chunk PPR partitions the input vector into $d$ chunks. However, this algorithm is designed for general inputs rather than specifically for sparse vectors. While the technique effectively reduces the computation time of PPR, our division method leverages the sparsity of the vector more efficiently.

In Section 4, we also provide a theoretical analysis of our algorithm, demonstrating that both the communication cost and execution time are independent of the vector's length. Specifically, given a privacy budget $\varepsilon$ and $d$ non-zero elements in the vector, we show that the communication cost is bounded by $O(\varepsilon d)$, while the computational complexity remains $O(d)$.

We emphasize that our communication cost depends only on the number of non-trivial elements, $d$, rather than the size of the vector, $n$. As a result, our cost is significantly lower than any previous algorithm designed for high-dimensional sparse vectors. Notably, our approach even outperforms the non-private case, which incurs a cost of $\Omega(d \log n)$.

Furthermore, since a smaller privacy budget corresponds to stronger privacy guarantees, our method achieves lower communication costs while ensuring better privacy protection. This contrasts with prior works on social network differential privacy [17, 10], where the communication cost increases as privacy improves. In those approaches, a smaller $\varepsilon$ leads to a larger number of non-trivial

elements in the obfuscated vector, thereby increasing the communication overhead.

We also emphasize that the expected computational cost of our algorithm scales linearly with the number of non-trivial elements, $d$. This represents a significant improvement over the mechanism proposed in [24], where the computation time grows exponentially with $d$.

From the compressed data, each element of the randomized response result can be retrieved in constant time.

Some applications of our algorithm are listed below. We also give the details of these application in Section 5.

*Graph/Social Network Data.* The private publication of social network adjacency lists is a fundamental tool in various graph privacy tasks. For instance, it serves as the first step in many synthetic graph generation frameworks; see [23] for a benchmark. Additionally, it plays a key role in subgraph counting, forming the initial phase of the general two-step mechanism [16], which was originally designed for triangle counting but has since been extended to cycle counting [11] and common neighbor estimation [9].

Given that this framework requires transmitting complete adjacency data to each user, reducing communication overhead is a crucial consideration. Our method allows users to publish their adjacency lists under edge-local differential privacy [27], with a communication cost proportional to their degree.

*Recommendation Systems.* A comprehensive survey on private recommendation systems can be found in [12], highlighting the growing interest in this field. In [2], the authors propose perturbing local user interaction ratings before submitting them to a central server that generates recommendations. However, their privacy model differs from ours, as it only protects the rating values while leaving their existence or absence unprotected. In [7], randomized response is employed for local obfuscation before transmitting user-item interactions to the central server.

We enable the publication of user-item interaction matrices, where the communication cost scales with the number of ratings.

*Genomic Data.* A common task in private genomics is the publication of single nucleotide polymorphisms (SNPs), which provide valuable insights into DNA variations and help assess genomic risk factors. Several methods have been developed for this purpose [15, 33, 34], though they rely on different privacy definitions than ours. In [32], randomized response is employed to publish SNP data under the $\varepsilon$-DP privacy model.

Our approach supports the release of single nucleotide polymorphisms (SNPs), with a communication cost proportional to the number of locations where the least common variation is present.

In Section 6, we present experimental results across all three application scenarios, highlighting the effectiveness and practicality of our mechanism. We verify that our communication and computation costs remain minimal even for vectors as large as 890,060 in size. Additionally, we confirm that a graph algorithm for triangle counting based on our framework achieves 100 times lower communication cost than any previous method, even in a small social network. With the same communication cost, our precision is $10^4$ better than any previous works.

## 2 PRELIMINARIES

### 2.1 Sampling with and without Replacement

In **sampling with replacement**, each draw is made independently from the entire population. If $n$ draws are performed, each with a probability of success $p$, the number of successes, denoted by $X$, follows a binomial distribution $\mathcal{B}(n, p)$. The probability of obtaining exactly $k$ successes is given by $\mathbb{P}(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$.

In contrast, **sampling without replacement** means that each draw is taken from the remaining unselected portion of the population. For a population of size $N$ containing $K$ successful elements, when drawing $n$ times, the number of successes, denoted by $Y$, follows a hypergeometric distribution $\mathsf{Hypergeometric}(N, K, n)$. The probability of observing exactly $k$ successes is given by $\mathbb{P}(Y = k) = \frac{\binom{K}{k}\binom{N-K}{n-k}}{\binom{N}{n}}$.

We will use the following theorem in our analysis.

**Theorem 1** (Theorem 4 of [13]). *Let $f$ be a continuous and convex function. If $Y \sim \mathsf{Hypergeometric}(N, K, n)$ and $X \sim \mathcal{B}(n, p)$ with $p = K/N$, then $\mathbb{E}[f(X)] \leq \mathbb{E}[f(Y)]$.*

Additionally, the standard formula for the moment generating function of $X \sim \mathcal{B}(n, p)$ is given as $\mathbb{E}\left[e^{tX}\right] = \left(1 - p + pe^t\right)^n$ (see Chapter 7.7 of [28]).

### 2.2 Counter-based Generators

Counter-based pseudo-random number generators (PRNGs) allow for the parallel generation of pseudo-random number sequences. Rather than generating all numbers from index 0 to $i - 1$ before obtaining the number at index $i$, each index can be generated independently without extra computational cost. This property is formally stated in Theorem 2.

**Theorem 2** (Counter-based PRNGs [29]). *There exist PRNGs $f$ such that, given a key $\mathcal{K}$ and an index $i$, the $i$-th element of the random sequence $f(\mathcal{K}, i)$ can be generated in constant time.*

There exist several efficient implementations of counter-based PRNGs (e.g., [29]), but in this article, we treat them as black-box mechanisms to simplify the presentation of our methods.

To formalize our notation, let $\mathcal{G}$ be an initialized counter-based PRNG. We denote by $\mathcal{G}^{(k)}$ the state of the generator after $k$ draws. Furthermore, for a given distribution $Q$, we define $\mathsf{Gen}(Q, \mathcal{G}^{(k)})$ as a sample drawn from $Q$ using the randomness generated by $\mathcal{G}$ at state $\mathcal{G}^{(k)}$.

### 2.3 Differential Privacy

**Definition 1** (Metric differential privacy [1]). Let d be a distance between datasets. A mechanism $\mathcal{M}$ satisfies $\varepsilon$-metric differential privacy if for any pair of datasets $D$ and $D'$, and for any possible outcome $S$ of $\mathcal{M}$, we have $\mathbb{P}\left[\mathcal{M}(D) \in S\right] \leq e^{\varepsilon d(D,D')}\mathbb{P}\left[\mathcal{M}(D') \in S\right]$.

The parameter $\varepsilon$ is called the privacy budget of $\mathcal{M}$.

Metric differential privacy was initially introduced as a convenient way to define privacy in metric spaces. However, it has also proven to be a highly general framework, as both the classic notions of local and central differential privacy can be seen as special cases of metric differential privacy under the appropriate choice of distance.

In this article, we consider an input space of $n$-dimensional vectors and define privacy using the Hamming distance $\mathcal{H}$. Specifically, for two vectors $v$ and $v'$, the distance $\mathcal{H}(v, v')$ is equal to the number of coordinates on which they differ.

The algorithm we will analyze in this article is randomized response, a method used to privately publish an entire vector of categorical data.

**Definition 2** (Randomized response [31]). For $\varepsilon > 0$ and a vector $v$ where each element of $v$, denoted by $v_1, \ldots, v_n$ belongs to $\{0, \ldots, k - 1\}$, the randomized response mechanism $\mathcal{R}$ with privacy budget $\varepsilon$ is defined as follows. For any $j \in \{0, \ldots, k - 1\}$, the probability of reporting $j$ instead of the true value $v_i$ is given by:

$$\mathbb{P}(\mathcal{R}(v_i) = j) = \begin{cases} \frac{e^\varepsilon}{e^\varepsilon + k - 1} & \text{if } v_i = j, \\ \frac{1}{e^\varepsilon + k - 1} & \text{otherwise.} \end{cases}$$

With the Hamming distance, randomized response satisfies $\varepsilon$-metric differential privacy.

### 2.4 Poisson Private Representation

In [24], the authors present a method called Poisson Private Representation (PPR) for converting any differentially private mechanism into a compressed version while preserving the original output distribution $P$. The transformed algorithm ensures an identical distribution to the original mechanism while having a communication cost of $O(\varepsilon)$.

This method leverages the shared random number generator results between the user and the central server, represented by draws $(Z_i)_{i \in \mathbb{N}}$ from a candidate distribution $Q$. The server should select $Q$ to closely approximate the true output distribution $P$ of the differentially private (DP) mechanism applied to the private data. While the server does not have access to the private data—and therefore cannot directly determine $P$—a practical approach is to use the output distribution of the same DP mechanism on arbitrary input data as an approximation.

With these shared draws, the user can transmit the index $K$ corresponding to the selected draw, which serves as the mechanism's output. The central server can then retrieve the output by computing $Z_K$. Notably, if counter-based PRNGs are used, this computation remains constant-time regardless of the value of $K$.

Theorem 3 establishes that this mechanism preserves both privacy and the original output distribution.

**Theorem 3** (Proposition 4.2 and Theorem 4.7 of [24]). *For an $\varepsilon$-metric differentially private mechanism $\mathcal{M}$, the PPR simulation of $\mathcal{M}$ with parameter $\alpha > 1$ satisfies $2\alpha\varepsilon$-metric differential privacy while ensuring that its output follows the same distribution as that of $\mathcal{M}$.*

Two types of communication occur during the protocol. The first is dedicated to establishing shared randomness between the user and the central server. This can be achieved, for instance, by transmitting a public instantiation key or using a predefined mechanism. Since this step can be initiated by either party and incurs only a small, constant communication cost, we exclude it from our communication analysis.

The second type of communication, which is our primary focus, involves the local user transmitting the value of $K$ to the central

server. This integer can be efficiently encoded using Huffman coding [14, 22], resulting in an expected communication cost of at most $\mathbb{E}[\log_2 K] + \log_2(\mathbb{E}[\log_2 K] + 1) + 2$ bits. Consequently, bounding the expected communication cost reduces to the problem of bounding $\mathbb{E}[\log_2 K]$.

**Theorem 4** (Theorem 4.3 of [24]). *For a mechanism with output distribution P, a candidate distribution Q, and a parameter $\alpha > 1$, the message K produced by PPR applied to P satisfies $\mathbb{E}\left[\log_2 K\right] \leq \mathrm{D}(P\|Q) + \frac{\log_2 3.56}{\min\{(\alpha-1)/2,1\}}$, where $\mathrm{D}(P\|Q)$ denotes the KL-divergence between P and Q.*

Another important consideration is the computational cost of the mechanism. On the server side, assuming the use of counter-based PRNGs, this cost remains constant and minimal. On the user side, however, the mechanism requires evaluating $\frac{\mathrm{d}P}{\mathrm{d}Q}(Z_i)$, which represents the ratio of the probability of obtaining $Z_i$ from $P$ to that from $Q$, for a large number of points $Z_i$ sampled from the probability distribution $Q$. To optimize this process, [24] introduces a reparameterization trick that employs a heap-based algorithm, ensuring the number of draws is bounded by $O\left(\sup_z \frac{\mathrm{d}P}{\mathrm{d}Q}(z)\right)$. Given that computing the probability ratio $\frac{\mathrm{d}P}{\mathrm{d}Q}(z)$ incurs a cost of $c$, the overall runtime of the PPR algorithm is $O\left(\sup_z \frac{\mathrm{d}P}{\mathrm{d}Q}(z) \cdot \left(c + \log \sup_z \frac{\mathrm{d}P}{\mathrm{d}Q}(z)\right)\right)$.

## 3 OUR ALGORITHM: COMPRESSION OF RANDOMIZED RESPONSE

The PPR mechanism can simulate any differentially private mechanism while reducing communication costs. Since randomized response is a differentially private mechanism, PPR can be directly applied to it. However, directly using PPR on the output of the randomized response method presents several challenges.

*Issue of PPR: Large Hamming Distance.* Since two distinct vectors of length $n$ can have a Hamming distance of up to $n$, the KL divergence between the candidate distribution and the actual distribution can reach $n\varepsilon$. Given this and considering Theorem 4, we can infer that the resulting communication cost scales as $O(n\varepsilon)$. Consequently, directly applying PPR to randomized response within this privacy model does not yield any improvement over naive randomized response.

The direct usage of PPR is challenging even when we focus on the common scenario where vectors are sparse or close to a reference. In this setting, the general structure of the vector is largely known before the user's data is published, with only a small number of coordinates differing from the reference. However, since the indices of these differing coordinates are unknown, the entire vector must still be published.

*Our Idea 1: Selection of the Candidate Distribution.* Let $d$ denote the number of differing coordinates from the reference. We select candidate distribution $Q$ as the randomized response applied to the reference vector. For instance, when the vector $v \in \{0, 1\}^n$ to be published is sparse, with only $d$ nonzero elements while the rest are zeros, the candidate distribution $Q$ is chosen as the distribution of the randomized response applied to the all-zero vector.

Under this choice, the KL divergence between the candidate distribution and the randomized response applied to the actual

vector is reduced to $d\varepsilon$. By carefully selecting the candidate distribution, we achieve a communication cost proportional to the vector's sparsity, similar to the non-private case.

*Issue of PPR: Computation Cost.* Since the number of draws (denoted as $\sup_z \frac{\mathrm{d}P}{\mathrm{d}Q}(z)$) in Section 2) increases exponentially with $d\varepsilon$, the computational cost of PPR also scales exponentially when applied to the randomized response result (see Section 8 of [24]). This makes it impractical for reasonable values of $d$. To reduce the number of required draws, we partition the adjacency vector into smaller groups. While the original PPR paper also employs chunking, their method cannot be directly applied here to ensure low computational cost. Their approach relies on contiguous chunks, but depending on the data structure, values may be concentrated in specific regions of the vector.

*Our Idea 2: Random Partitioning Strategy.* To reduce the number of draws, we divides the vector into smaller chunks using random partitioning strategy, ensuring that each chunk has a small degree on average. PPR is then applied independently to each chunk. By carefully selecting the chunk size, we achieve a communication cost proportional to the degree while maintaining a computational cost that also scales with the degree.

---

**1** **Function** EncodeRR
    **Input:** 1) A list of indices corresponding to the non-trivial entries of the input vector, denoted as $(x_1, \ldots, x_d)$; 2) the values of the input vector at these indexes $(v_1, \ldots, v_d)$; 3) The reference vector values at these indexes $(c_1, \ldots, c_d)$; 4) The number of chunks $m$; 5) A privacy budget $\varepsilon$; 6) A public random permutation function $\varphi$
    **Output:** The encoded result $(K_1, \ldots, K_m)$
**2**    $y_1, \ldots, y_d \leftarrow \varphi(x_1), \ldots, \varphi(x_d)$;
**3**    $s \leftarrow \lceil n/m \rceil$;
**4**    For each $i \in [1, d]$, compute the Euclidean division of $y_i$ by $s$, yielding $(q_i, r_i)$;
**5**    **for** $j \leftarrow 1$ **to** $m$ **do**
**6**        $S_j \leftarrow \{(r_i, u_i, c_i) \mid q_i = j\}$;
**7**        $K_j \leftarrow \mathrm{PPR}(S_j, \varepsilon, \alpha)$;
**8**    **end**
**9**    **return** $(K_1, \ldots, K_m)$

**Algorithm 1:** Encodes the compressed randomized response of a vector

---

Our encoding algorithm is given in Algorithm 1. It begins by randomly permuting the vector using the function $\varphi$, and divide it into $m$ chunks. The non-trivial elements of those $m$ chunks are denoted by the set $S_1, \ldots, S_m$. Next, we apply the PPR method [24] to each chunk and return the resulting list as the output.

Our decoding algorithm is presented in Algorithm 2. To determine the value of the input vector at index $i$, we first compute the chunk number $q$ and the position of $i$ within the chunk, denoted by $r$, using the public function $\varphi$ and $s$. We then perform the decoding using the same approach as PPR.

---

**1 Function** DecodeRR
    **Input:** A list of compressed indexes $(K_1, \ldots, K_m)$, a list
        of distributions $(Q_1, \ldots, Q_n)$, the index $i$ that
        one wants to access
    **Output:** The value of the vector at index $i$
**2**    $j \leftarrow \varphi(i)$;
**3**    Let $(q, r)$ be the result of the Euclidean division of $j$ by $s$;
**4**    **return** $\mathsf{Gen}(Q_i, \mathcal{G}^{(sK_q + r)})$

---

**Algorithm 2:** Decodes the compressed randomized response of a vector

*Our Idea 3: Efficient Calculation of* $\frac{dP}{dQ}(Z_i)$. As noted in Section 2, PPR requires multiple evaluations of $\frac{dP}{dQ}(Z_i)$. A naive approach to computing $\frac{dP}{dQ}(Z_i)$ involves generating the entire randomized vector and comparing its probability under both distributions. Even when Algorithm 1 reduces the vector size from $n$ to $s$, $s$ typically remains of the same order as $n$, making this method computationally inefficient and leading to significant computation time.

In Algorithm 3, we propose an efficient approach leverages the fact that this probability ratio depends only on the values of the draw at the indices where the private vector differs from the reference.

Let $Z_i = (z_1, \ldots, z_s)$ be a sparse vector, and $x_1, \ldots, x_{d'}$ are the indices on which the input vector has non-trivial values for all $1 \leq j \leq d'$. By the independence of the randomized response mechanism, we observe that $\mathbb{P}_P[Z_i]/\mathbb{P}_Q[Z_i] = \prod_{j=1}^{d'} \mathbb{P}_{P_j}[Z_{x_j}]/\mathbb{P}_{Q_j}[Z_{x_j}]$, where $P_j$ represents the probability distribution of the randomized response result derived from the $x_j$-th element of the input vector, and $Q_j$ corresponds to the one obtained from the reference vector. Since $z_j$ can be generated independently and is not required for the calculation when $j \notin \{x_1, \ldots, x_{d'}\}$, we can omit $z_j$ in such cases. Leveraging this observation, we can bypass generating the full vector and instead compute the ratio using only these $d'$ specific coordinates. The calculation time $c$ is $O(d')$.

---

**1 Function** ProbabilityRatio
    **Input:** 1) A list of indices corresponding to the
        non-trivial entries of the input vector, denoted as
        $(x_1, \ldots, x_{d'})$; 2) The values of the input vector at
        these indices, represented as $(v_1, \ldots, v_{d'})$; 3) The
        reference vector values at these indices, given by
        $(c_1, \ldots, c_{d'})$; 4) The value of $Z_i$ at these indices,
        given by $(z_1, \ldots, z_{d'})$; 5) A privacy budget $\varepsilon$.
    **Output:** $\frac{dP}{dQ}(Z_i)$
**2**    $ratio \leftarrow 1$;
**3**    **for** $i \leftarrow 1$ **to** $d'$ **do**
**4**       **if** $z_i = v_i$ **then** $ratio \leftarrow ratio \times e^{\varepsilon}$ ;
**5**       **if** $z_i = c_i$ **then** $ratio \leftarrow ratio \times e^{-\varepsilon}$ ;
**6**    **end**
**7**    **return** $ratio$

---

**Algorithm 3:** Calculate the ratio $dP/dQ$ at a given state $Z_i$

## 4 THEORETICAL ANALYSIS

First, we give the privacy result for our algorithm.

**Theorem 5.** *Algorithm 1 satisfies $2\alpha\varepsilon$-metric differential privacy.*

PROOF. Each independent PPR satisfies $2\alpha\varepsilon$-metric differential privacy, as stated in Theorem 3. Since the mechanism partitions the indices into $m$ groups, its overall privacy guarantee follows from the parallel composition property of differential privacy [26, 25]. □

The number of chunks, $m$, is a tunable parameter that balances communication cost and execution time. Our analysis focuses on the case where $m = \beta\varepsilon d$, with $\beta$ as a parameter. We demonstrate that under this setting, the communication cost is of order $O(\varepsilon d)$, while the computation cost remains $O(d)$.

**Theorem 6.** *The communication cost of our algorithm is $O(\varepsilon d)$, where $d$ represents the number of indices where the input vector differs from the reference vector.*

**Lemma 1.** *For $\mathcal{M}$ a metric differential private mechanism, and $P$ and $Q$ two distributions resulting from $\mathcal{M}$ applied on two datasets at distance $d$, then $\sup_z \frac{dP}{dQ}(z) \leq e^{\varepsilon d}$ and $D(P \| Q) \leq \varepsilon d$.*

PROOF. Since $P$ and $Q$ are derived from the same metric differentially private mechanism applied to two datasets separated by a distance $d$, it follows that $\frac{dP}{dQ}(z) \leq e^{\varepsilon d}$, for all $z$. Using this result, we obtain $D(P \| Q) = \mathbb{E}_{Z \sim P}\left[\log\left(\frac{dP}{dQ}(Z)\right)\right] \leq \varepsilon d$. □

PROOF OF THEOREM 6. The random public permutation function can be obtained through shared common knowledge. Additionally we suppose that all users (including the central server) are aware of the reference vector. This leads us to only consider the sharing of $(K_1, \ldots, K_m)$ for the communication cost.

We define $d_i = |S_i|$, ensuring that $d_1 + \cdots + d_m = d$. Each $d_i$ represents the distance from the reference point of the $i$-th vector. PPR is applied independently to every chunk $S_i$. Thus, each $K_i$ verifies $\mathbb{E}\left[\log_2 K_i\right] \leq D(P_i \| Q_i) + \frac{\log_2 3.56}{\min\{(\alpha - 1)/2, 1\}}$ and $D(P_i \| Q_i) \leq \varepsilon d_i$ using Lemma 1. This gives a total expected communication cost of $\sum_{i=1}^{m} \left[\mathbb{E}\left[\log_2 K_i\right] + \log_2(\mathbb{E}\left[\log_2 K_i\right] + 1) + 2\right]$
$\leq \sum_{i=1}^{m} \left[D(P_i \| Q_i) + O(1) + \log_2\left(D(P_i \| Q_i) + O(1)\right)\right]$
$\leq \sum_{i=1}^{m} \left[\varepsilon d_i + O(1) + \log_2\left(\varepsilon d_i + O(1)\right)\right] = O(\varepsilon d)$. □

**Theorem 7.** *The computational cost of the compressed randomized response proposed in this work is $O(d)$, where $d$ represents the number of indices where the input vector differs from the reference vector.*

PROOF. The primary contributor to computational cost is the runtime of PPR for the various $S_i$. Therefore, our analysis will focus on these computations. For any $i \in [1, m]$, we first observe that the complexity $c_i$ of Algorithm 3 applied to $S_i$ is $\Theta(d_i)$. Consequently, the expected computational cost $C_i$ of running PPR depends on $d_i$ and is given by $C_i(d_i) = O\left(e^{\varepsilon d_i}(d_i + \varepsilon d_i)\right) = O\left(d_i e^{\varepsilon d_i}\right)$.

We need to analyze the distribution of $d_i$ and its impact on the expected value of $C_i$. The indices of the vector are randomly shuffled before being divided into $m$ groups. In the worst-case scenario, the subset $S_i$ forms a contiguous block of size $\lceil n/m \rceil$. Under this scenario, $d_i$ follows a hypergeometric distribution with parameters

$n, d, \lceil n/m \rceil$. Furthermore, since $C_i$ is a convex function, we can apply Theorem 1. This result establishes that the expected value of $C_i(d_i)$ is upper-bounded by the expected value of $C_i(Y)$, where $Y \sim \mathcal{B}(N, p)$ with $N = \lceil n/m \rceil$ and $p = d/n$.

We define the function $f(x) = xe^{\varepsilon x}$ and compute the expected value of $f(Y)$. This allows us to bound the complexity of Algorithm 1 by $O\left(m \cdot \mathbb{E}[f(Y)]\right)$, where $\mathbb{E}[f(Y)] = \sum_{i=0}^{N} ie^{\varepsilon i}p^i(1 - p)^{N-i}\binom{N}{i}$. This expression can be viewed as a function of $e^{\varepsilon}$, which we denote as $g$. Additionally, we introduce the function $h(x) = \sum_{i=0}^{N} x^i p^i (1-p)^{N-i}\binom{N}{i}$. From the moment-generating function formula for the binomial distribution, we obtain $h(x) = (1+(x-1)p)^N$. Differentiating $h(x)$, we find $h'(x) = pN(1 + (x - 1)p)^{N-1}$. Since $g(x)$ satisfies the relation $xh'(x) = g(x)$, we conclude that:

$$
\mathbb{E}\left[f(Y)\right] = e^{\varepsilon}pN \left(1 + (e^{\varepsilon} - 1)p\right)^{N-1}
$$
$$
\leq \frac{d}{n}\left(\frac{n}{m} + 1\right)\left(1 + (e^{\varepsilon} - 1)\frac{d}{n}\right)^{\lceil \frac{n}{m} \rceil - 1}
$$
$$
\leq \left(\frac{1}{\beta\varepsilon} + \frac{d}{n}\right)\exp\left[\frac{n}{m}\ln\left(1 + (e^{\varepsilon} - 1)\frac{d}{n}\right)\right]
$$
$$
\leq \left(\frac{1}{\beta\varepsilon} + \frac{d}{n}\right)\exp\left[\frac{n}{m}(e^{\varepsilon} - 1)\frac{d}{n}\right]
$$
$$
= \left(\frac{1}{\beta\varepsilon} + \frac{d}{n}\right)\exp\left(\frac{e^{\varepsilon} - 1}{\beta\varepsilon}\right) = O\left(\frac{1}{\beta\varepsilon}\right)
$$

The final step of the derivation follows from the fact that $\exp\left(\frac{e^{\varepsilon}-1}{\beta\varepsilon}\right) = O(1)$ when $\beta \geq 1$ and $\varepsilon$ is close to zero. Since Algorithm 1 executes a total of $m = \beta\varepsilon d$ instances of PPR, the overall computational complexity of the algorithm is given by $O\left(m \cdot \mathbb{E}[f(Y)]\right)$, which simplifies to $O\left(d\right)$. □

## 5 POTENTIAL APPLICATIONS

In this section, we will explore various scenarios in which our algorithm can be applied. Randomized response serves as a fundamental component of differentially private algorithms. Consequently, our mechanism is particularly useful in private applications involving large volumes of data where communication costs are a concern. The only prerequisites are (1) the existence of a reference vector and (2) the adoption of metric differential privacy as the privacy framework. Regarding the first requirement, in most real-world scenarios, the server typically has some prior knowledge of the information held by the user. In the following discussion, we will examine examples of how this prior knowledge about the secret vector can be transformed into a reference vector.

*Graph/Social Network Information.* The first scenario we examine is the publication of adjacency lists. In this setting, users seek to privately share their list of neighbors in a graph. The primary framework used to protect such information is edge-local differential privacy [27], which aligns with metric differential privacy by defining distance as the Hamming distance between two adjacency vectors. Since most real-world graphs' adjacency vectors are sparse, it is generally known in advance that the majority of bits in the adjacency vector will be zeros. Therefore, we adopt a reference vector consisting entirely of zeros for our mechanism. According to Theorem 6, this choice leads to a communication cost

proportional to $d$, the degree of the node—typically around a thousand—rather than $n$, the total number of nodes in the graph, which usually reaches several million, as would be required for a naive randomized response.

*Recommendation System.* In this scenario, each user holds a set of ratings for certain items. Those ratings forms a vector called user-item interaction. Specifically, for each item in the set of possible items, a user has either not provided a rating or has assigned a score from a finite set of values. Our goal is to publish these ratings while preserving metric differential privacy, where the distance between two vectors is defined as the number of items for which the ratings differ or are present in only one of the two vectors.

Since the total number of possible items, represented by the size of the user-item interaction vector $n$, is typically very large, while each user has rated only a small subset, we use an empty rating vector (where no items have been rated) as the reference. Consequently, the number of non-trivial elements in the input vectors, denoted as $d$, corresponds to the number of items a user has rated, which is significantly smaller than $n$. This allows our algorithm to achieve a communication cost proportional to the number of rated items rather than the total number of possible items, substantially reducing overhead.

*Genomic Information.* We focus on the publication of single-nucleotide polymorphisms (SNPs), which represent variations of a single nucleotide in the genome. In this setting, both the central server and users have access to a list of possible SNP locations, along with the most common nucleotide variation at each location. The reference vector is constructed using these most frequent variations, while the input vector represents each user's specific SNP data. Consequently, the vector size $n$ corresponds to the total number of locations, while $d$ represents the number of locations where a user's genomic information differs from the most frequent variation. It is known that $d \ll n$ in this type of dataset.

## 6 EXPERIMENTAL RESULTS

In this section, we conduct experiments on the three applications identified in the previous section. The code for these experiments is available in the following repository: https://anonymous.4open.science/r/Metric-DP-Compression-8362. All timed experiments are conducted on a MacBook Pro (14-inch, 2021) equipped with an M1 Pro chip featuring an 8-core CPU and 32GB of memory. We do not include a comparison with PPR or chunk PPR from [24], as their execution time is prohibitively high across all experimental settings. Additionally, we do not compare the communication cost and execution time with the randomized response technique, as its communication cost is significantly higher than our method and can be theoretically predicted. Moreover, its execution time remains consistently low since it only involves bit flipping.

### 6.1 Recommendation Systems

We conduct our experiments using the MovieLens 32M dataset [8], which contains 32 million ratings for 87,585 movies from 200,948 users. For all experiments, we randomly select 1,000 users from this dataset and apply our algorithm to their rating lists. The vector size (representing the number of movies), $n$, is 87,585, while the

number of non-trivial elements (representing the number of ratings per user), $d$, ranges from a few to 3,500.

*Upload Cost and Execution Time.* First, in Figure 1, we present the communication cost required for users to upload their randomized response to the server, along with the execution time, using the default parameters: $\varepsilon = 1$, $\alpha = 2$, and $\beta = 2$.

The results indicate that the upload cost scales linearly with the number of ratings, as expected, with a slope of approximately 2. The results also confirm that the upload cost of our algorithm is even lower than that of non-private publication of the adjacency vector when using the adjacency list format. While the execution time also increases with the number of ratings, the high variance reduces the strength of this correlation. We have verified that the high variance is not due to the variance of $d_i$ in our algorithm but rather stems from the PPR. Despite the large variance, all 1,000 executions maintain a manageable execution time.



**Figure 1: Communication cost and the execution time incurred by our algorithm for 1000 users of the Movie Lens dataset as a function of their number of movie rated**

*Results on Different Privacy Budget $\varepsilon$.* In this experiment, all parameters remain at their default values except $\varepsilon$. The results of this experiment are presented in Figure 2 for the upload cost and Figure 3 for the execution time. These results indicate that as the privacy budget increases, both the upload cost and execution time of the algorithm increase, while the variance of the execution time decreases.
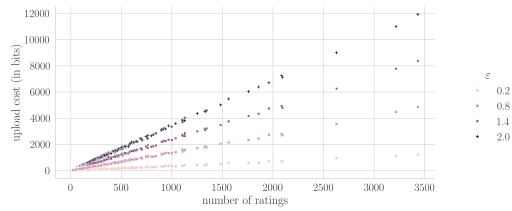


**Figure 2: Communication cost of our algorithm for 1,000 users in the MovieLens dataset across different privacy budget values**

*Results on Different Values of Parameter $\beta$.* Recall that the number of chunks, $m$, is defined as $\beta \varepsilon d$. In the appendix, we present experimental results for different values of $\beta$, which confirm that setting the default value to 2 is a reasonable choice.
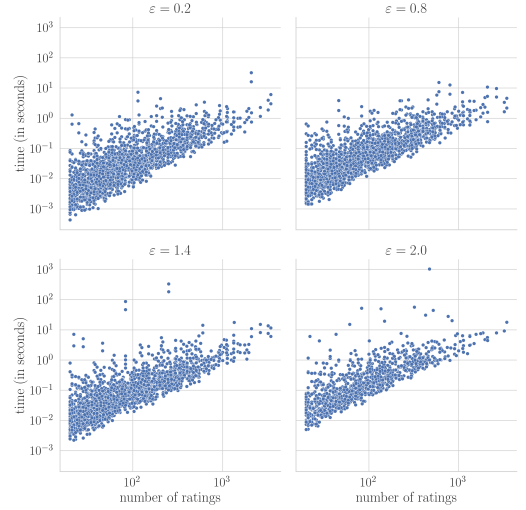


**Figure 3: Execution time of our algorithm on 1000 users of the Movie Lens dataset for different values of the privacy budget**

*Estimation of the Number of Common Items.* To prove the accuracy of our method, we also evaluate it on the task of computing the number of common items between 2 users. To this end, we randomly select pairs of users in MovieLens and estimate their number of items in common with classic randomized response and our algorithm, which is called compressed RR in the figure.
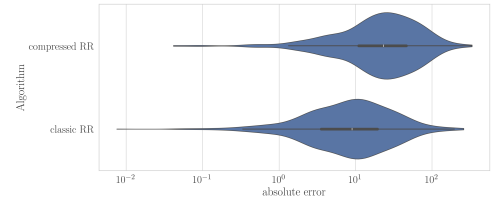


**Figure 4: The absolute error in estimating the number of common neighbors across 1,000 user pairs in the MovieLens dataset**

The results in Figure 4 show that accuracy experiences only a slight decline when using our algorithm, which is expected. Our method can achieve the same results as randomized response but with a larger privacy budget. Consequently, when the budget is fixed, the accuracy is slightly reduced.

## 6.2 Genomic Data

The second application we examine is the publication of SNPs by users. For our experiments, we use chromosome 22 data from the Phase 3 release of the 1000 Genomes Project [3].

This dataset consists of 1,064,502 locations. While real SNP information is not available due to its sensitivity, we have probability values indicating the likelihood of a user having a variation from the most frequent nucleotide at each location. We use these probabilities to generate synthetic user data. All probabilities are greater

than $10^{-4}$. Our method is particularly beneficial for sparse datasets, so we exclude locations where variations are too frequent from our experiments. This exclusion does not pose a limitation, as the principle of parallel composition allows us to publish frequent variations using classical randomized response while applying our algorithm to the rarer ones—without splitting the privacy budget. As a result, we retain only variations with a frequency below $10^{-2}$, leading to $n = 890,060$ remaining variations, which represents over 83% of the dataset.

In our experiments, we generate 1,000 synthetic users and apply our algorithm to their SNP lists, with the number of variations ranging from 800 to 1,100.

In Figure 5, we plot the resulting upload cost and execution time as functions of the number of variations each user possesses. Since the number of variations falls within a narrower range, the distribution is more clearly visible compared to the recommendation system. Notably, the variance in execution time is too high to reveal a clear trend. However, the upload cost remains within a 10% range of its average value, indicating stable performance. In all of the results, our algorithm exhibits efficient performances both in the upload cost and the execution time.
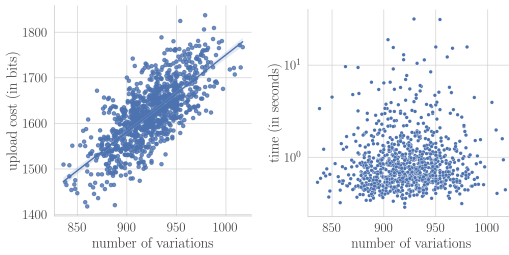


**Figure 5: The communication cost and the execution time of our algorithm for 1000 randomly generated SNPs sequences as a function of the number of variations**

### 6.3 Social Networks

*Upload Cost and Execution Time.* In Appendix, we present our upload cost and execution time on the Google+ dataset, which consists of 107,614 nodes. The results align closely with those observed in the recommendation system and genomic data experiments.

*Triangle Counting.* For the experiment on triangle counting we chose to conduct them on the Wikipedia graph [20, 19]. This graph contains $n = 7,115$ nodes and 103,689 edges. We are unable to conduct this experiment on the Google+ dataset because the randomized response technique requires excessive memory, making it infeasible to run the algorithm within our computational environment.

We use the two-step mechanisms described in [16] to privately estimate the number of triangles in the graph, with one key modification: we replace the classical randomized response with our algorithm. In this two-step mechanism, all users must download the randomized response results of every other user to their local storage. As a result, most of the communication cost comes from these download costs. Therefore, unlike other experiments

where we compare upload costs, we focus on download costs in this evaluation.

For comparison, we evaluate our method against ARROne [17] and GroupRR with CSS [10]. These algorithms include a sampling parameter that balances communication cost and accuracy. To demonstrate the full range of their capabilities, we compute their $\ell_2$-errors for various values of this parameter.

In contrast, our method does not require such a parameter, so we represent only a single point in the results, corresponding to the average download cost and the average $\ell_2$-error, defined as the square root of the sum of squared errors over 10 runs. The results are presented in Figure 6.
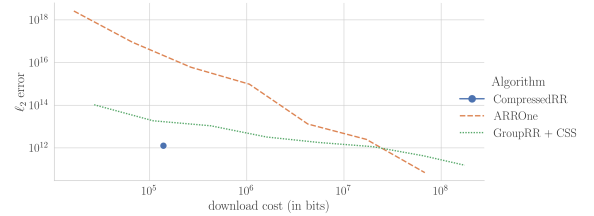


**Figure 6: The average $\ell_2$ error in the estimation of the number of triangles in the Wiki graph for 3 different algorithms**

We observe that, for the same level of accuracy, our method reduces communication cost by a factor of more than 100. Furthermore, at the communication cost used by our method, the error of GroupRR is over 10 times higher, while ARR exhibits an error more than $10^4$ times worse. We expect these improvements to be even more pronounced for larger or sparser graphs.

## 7 CONCLUSION

Randomized response is one of the most widely used algorithms for protecting users' sensitive information under metric and local differential privacy, with numerous potential applications. However, for sparse vectors, this method is inefficient in both communication and storage costs.

Although several studies have proposed solutions to mitigate this issue [17, 10], we fully resolve it by achieving an even lower cost than non-private communication. This significantly expands the applicability of randomized response. Our algorithm is built on an information-theoretic approach inspired by PPR. While PPR is known to compress information published under differential privacy, its compression rate is approximately $1/\epsilon$, where $\epsilon$ is the privacy budget. In contrast, our compression rate is $n/(\epsilon d)$, where $n$ is the vector size and $d$ is the number of non-trivial values in the sparse vector. This rate is significantly higher than that of PPR.

Our algorithm is the first to demonstrate how an information-theoretic approach can drastically reduce communication costs in differential privacy applications.

# REFERENCES

[1] Miguel E Andrés, Nicolás E Bordenabe, Konstantinos Chatzikokolakis, and Catuscia Palamidessi. 2013. Geo-indistinguishability: differential privacy for location-based systems. In *SIGSAC 2013*, 901–914.

[2] Arnaud Berlioz, Arik Friedman, Mohamed Ali Kaafar, Roksana Boreli, and Shlomo Berkovsky. 2015. Applying differential privacy to matrix factorization. In *RecSys 2015*, 107–114.

[3] 1000 Genomes Project Consortium et al. 2015. A global reference for human genetic variation. *Nature*, 526, 7571, 68.

[4] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam D. Smith. 2006. Calibrating noise to sensitivity in private data analysis. In *TCC 2006*, 265–284.

[5] Ulfar Erlingsson, Vasyl Pihur, and Aleksandra Korolova. 2014. RAPPOR: Randomized aggregatable privacy-preserving ordinal response. In *SIGSAC 2014*, 1054–1067.

[6] Vitaly Feldman and Kunal Talwar. 2021. Lossless compression of efficient private local randomizers. In *ICML 2021*, 3208–3219.

[7] Chen Gao, Chao Huang, Dongsheng Lin, Depeng Jin, and Yong Li. 2020. DPLCF: Differentially private local collaborative filtering. In *SIGIR 2020*, 961–970.

[8] F. Maxwell Harper and Joseph A. Konstan. 2015. The MovieLens datasets: History and context. *ACM Trans. Interact. Intell. Syst.*, 5, 4.

[9] Yizhang He, Kai Wang, Wenjie Zhang, Xuemin Lin, and Ying Zhang. 2024. Common neighborhood estimation over bipartite graphs under local differential privacy. *PACMMOD 2024*, 2, 6, 1–26.

[10] Quentin Hillebrand, Vorapong Suppakitpaisarn, and Tetsuo Shibuya. 2023. Communication cost reduction for subgraph counting under local differential privacy via hash functions. *arXiv preprint arXiv:2312.07055*.

[11] Quentin Hillebrand, Vorapong Suppakitpaisarn, and Tetsuo Shibuya. 2025. Cycle counting under local differential privacy for degeneracy-bounded graphs. *STACS 2025*.

[12] Yassine Himeur, Shahab Saquib Sohail, Faycal Bensaali, Abbes Amira, and Mamoun Alazab. 2022. Latest trends of security and privacy in recommender systems: a comprehensive review and future perspectives. *Computers & Security*, 118, 102746.

[13] Wassily Hoeffding. 1963. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58, 301, 13–30.

[14] David A Huffman. 1952. A method for the construction of minimum-redundancy codes. *Proceedings of the IRE*, 40, 9, 1098–1101.

[15] Mathias Humbert, Erman Ayday, Jean-Pierre Hubaux, and Amalio Telenti. 2014. Reconciling utility with privacy in genomics. In *WPES 2014*, 11–20.

[16] Jacob Imola, Takao Murakami, and Kamalika Chaudhuri. 2021. Locally differentially private analysis of graph statistics. In *USENIX Security 2021*, 983–1000.

[17] Jacob Imola, Takao Murakami, and Kamalika Chaudhuri. 2022. Communication-efficient triangle counting under local differential privacy. In *USENIX Security 2022*, 537–554.

[18] Shiva Prasad Kasiviswanathan, Homin K. Lee, Kobbi Nissim, Sofya Raskhodnikova, and Adam Smith. 2011. What can we learn privately? *SIAM Journal on Computing*, 40, 3, 793–826.

[19] Jure Leskovec, Daniel Huttenlocher, and Jon Kleinberg. 2010. Predicting positive and negative links in online social networks. In *WWW 2010*, 641–650.

[20] Jure Leskovec, Daniel Huttenlocher, and Jon Kleinberg. 2010. Signed networks in social media. In *CHI 2010*, 1361–1370.

[21] Jure Leskovec and Julian Mcauley. 2012. Learning to discover social circles in ego networks. *NIPS 2012*, 25, 1–9.

[22] Cheuk Ting Li and Abbas El Gamal. 2018. Strong functional representation lemma and applications to coding theorems. *IEEE Transactions on Information Theory*, 64, 11, 6967–6978.

[23] Shang Liu, Hao Du, Yang Cao, Bo Yan, Jinfei Liu, and Masatoshi Yoshikawa. 2025. PGB: Benchmarking differentially private synthetic graph generation algorithms. *ICDE 2025*.

[24] Yanxiao Liu, Wei-Ning Chen, Ayfer Özgür, and Cheuk Ting Li. 2024. Universal exact compression of differentially private mechanisms. *NeurIPS 2024*.

[25] Pasin Manurangsi and Warut Suksompong. 2023. Differentially private fair division. *AAAI 2023*, 5814–5822.

[26] Frank D McSherry. 2009. Privacy integrated queries: an extensible platform for privacy-preserving data analysis. In *SIGMOD 2009*, 19–30.

[27] Zhan Qin, Ting Yu, Yin Yang, Issa Khalil, Xiaokui Xiao, and Kui Ren. 2017. Generating synthetic decentralized social graphs with local differential privacy. In *CCS 2017*, 425–438.

[28] Sheldon M Ross. 1976. *A first course in probability*. Vol. 2. Macmillan New York.

[29] John K Salmon, Mark A Moraes, Ron O Dror, and David E Shaw. 2011. Parallel random numbers: as easy as 1, 2, 3. In *SC 2011*, 1–12.

[30] Abhin Shah, Wei-Ning Chen, Johannes Balle, Peter Kairouz, and Lucas Theis. 2022. Optimal compression of locally differentially private mechanisms. In *AISTATS 2022*, 7680–7723.

[31] Yue Wang, Xintao Wu, and Donghui Hu. 2016. Using randomized response for differential privacy preserving data collection. In *EDBT/ICDT 2016* number 35.

[32] Akito Yamamoto and Tetsuo Shibuya. 2024. Privacy-Optimized Randomized Response for Sharing Multi-Attribute Data. In *ISCC 2024*, 1–8.

[33] Emre Yilmaz, Erman Ayday, Tianxi Ji, and Pan Li. 2020. Preserving genomic privacy via selective sharing. In *WPES 2020*, 163–179.

[34] Emre Yilmaz, Tianxi Ji, Erman Ayday, and Pan Li. 2022. Genomic data sharing under dependent local differential privacy. In *CODASPY 2022*, 77–88.

## APPENDIX: ADDITIONAL EXPERIMENTAL RESULTS

### Results on Movie Lens Dataset for Different Values of $\beta$

In this experiment, all parameters remain at their default values except for $\beta$. The results are presented in Figure 7 for the upload cost and Figure 8 for the execution cost. The findings indicate that increasing $\beta$ leads to a higher upload cost. For execution time, both its value and variance decrease as $\beta$ increases. However, this trend is not observed between $\beta = 2$ and $\beta = 4$, which led us to select $\beta = 2$ as the default value.



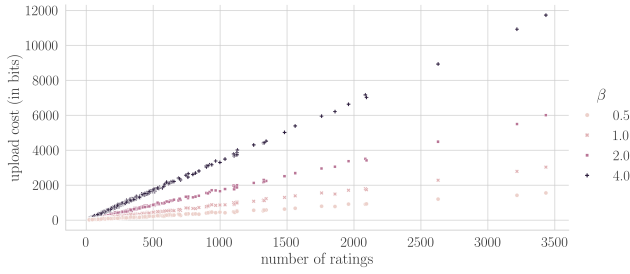**Figure 7: The communication cost of our algorithm on 1000 users of the Movie Lens dataset for different values of the parameter $\beta$.**
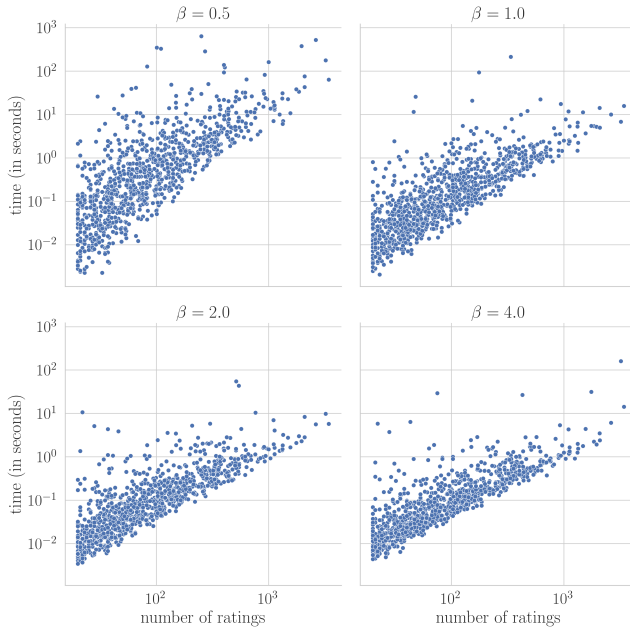


**Figure 8: The execution time of our algorithm on 1000 users of the Movie Lens dataset for different values of the parameter $\beta$.**

### Experiments on Google+ Dataset

We conducted our experiments on the Google+ dataset [21], where nodes represent users and edges indicate connections between users within a circle. The resulting graph consists of $n = 107,614$ nodes and 13,673,453 edges. The degree $d$ ranges from 1 to 5,000.

To publish the complete adjacency matrix of an unordered graph, it is sufficient for each user to disclose only their connections with nodes having smaller indices than their own [10]. Based on this principle, we applied our algorithm to the adjacency vector, retaining only the 1s corresponding to connections with lower-indexed nodes. Figure 9 presents the results for 1,000 randomly selected users from the graph.
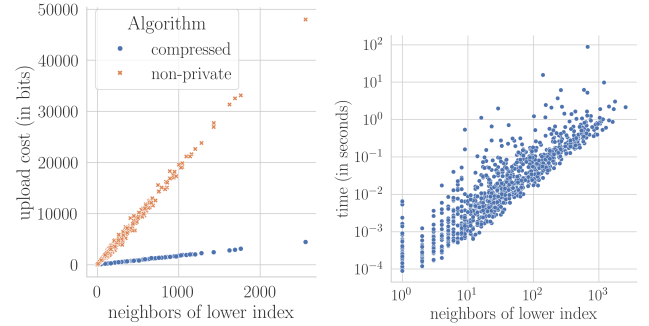


**Figure 9: The communication cost and the execution time incurred by our algorithm for 1,000 users of the Google+ dataset as a function of their number of neighbors of smaller index**

Similar to user-item interactions in recommendation systems, we observe that both the upload cost and execution time increase with the number of neighbors having smaller indices. Additionally, the variance in execution time is higher than in upload cost. However, we also note that both the upload cost and execution time remain low across all data points.