

Aggregation and Transformation of Vector-Valued Messages in the Shuffle Model of Differential Privacy

Mary Scott, Graham Cormode and Carsten Maple (IEEE Transactions on Information Forensics and Security)

Differential privacy (DP) is a technique that provides a rigorous and provable privacy guarantee for aggregation and release. The Shuffle Model for DP has been introduced to balance the accuracy of local-DP algorithms with the privacy risks of central-DP. In this work we firstly provide a single message protocol for the summation of real vectors in the Shuffle Model. We then improve this bound through the implementation of a Discrete Fourier Transform, greatly minimizing the initial error at the expense of a small loss in accuracy through the transformation itself. This work will further the exploration of more sophisticated structures such as matrices and higher-dimensional tensors in this context, both of which are reliant on the functionality of the vector case.

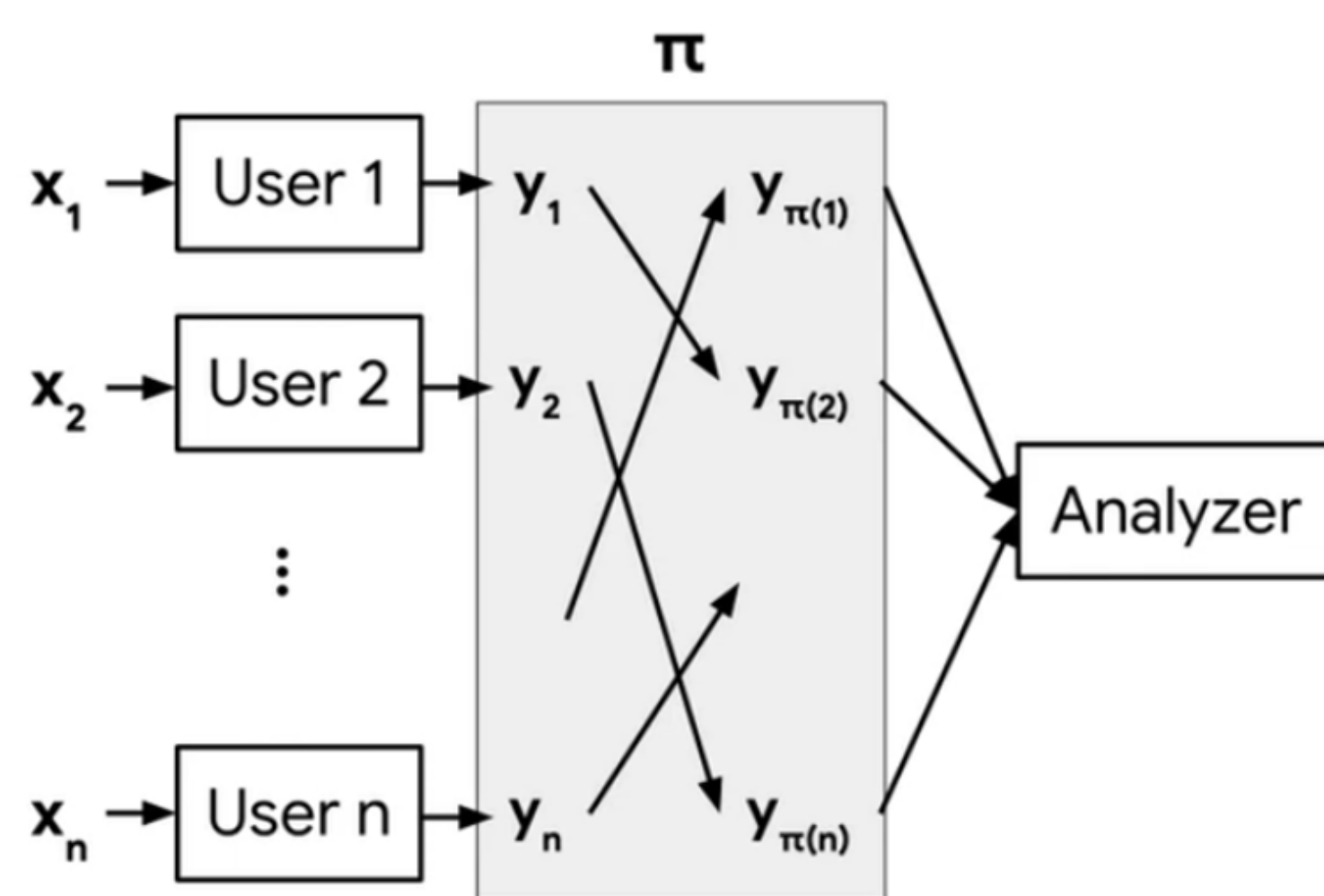


Figure 1: The addition of a shuffling step to the Local Model of DP (represented by the permutation π).

Introduction

One of the fundamental challenges of data analysis is the careful balance of acquiring as much utility from a dataset as possible, whilst simultaneously providing a strong guarantee of privacy to each individual affected.

A function applied to a dataset is *differentially private* (DP) if, with the removal of any individual from the dataset, the output of the function does not change more than a small multiplicative factor ϵ .

Definition [1]: A randomized function M gives ϵ -differential privacy if for all datasets x and y differing on at most one element, and all $S \subseteq \text{Range}(M)$,

$$\Pr[M(x) \in S] \leq \exp(\epsilon) \Pr[M(y) \in S].$$

Related Work

The majority of research in the field of DP has focused on two contrasting models:

- **Centralized Model:** users submit their sensitive personal information directly to a *trusted* central data collector, who adds *random noise* to the raw data to provide DP, before assembling and analyzing the aggregated results.
- **Local Model:** DP is guaranteed when each user applies a *local randomizer* to add random noise to their data before it is submitted. There is no need for a trusted party, but the level of noise required per user for the same privacy guarantee is much higher.

In recent years researchers have tried to create intermediate models that reap the benefits of both. In 2019, Cheu *et al.* [2] formalized the Shuffle Model to connect an additional shuffle step to the Local Model. This step is shown by the grey rectangle in Figure 1.

Vector Sum in the Shuffle Model

Our first contribution is a new protocol in the Shuffle Model for the private summation of vector-valued messages. This protocol extends an existing result from Balle *et al.* [3] by permitting the n users to each submit a vector of real numbers rather than being restricted to submitting a scalar.

The local randomizer applies a generalized *randomized response* mechanism that:

- returns the true message x_i with probability $1 - \gamma$,
- and a uniformly random message with probability γ .

It is necessary to find an appropriate γ to optimize the proportion of random messages that are submitted, and also guarantee DP. The resulting estimator is unbiased and has normalized mean squared error (MSE) $O_{\epsilon, \delta}(d^{8/3}n^{-5/3})$, where d is the dimension of each vector.

Transforming Summation in the Shuffle Model

An orthonormal transformation can be used to further tighten the bound we have obtained for private summation. The Discrete Fourier Transform (DFT) concentrates information about signals with a particular property into a small number of coefficients.

Our second contribution, which we call the Fourier Summation Algorithm (FSA), combines the private summation protocol with the DFT from Rastogi and Nath in the centralized case [4]. This improves the accuracy of the tight bound to $O_{\epsilon, \delta}(m^{8/3}n^{-5/3})$, where m represents the number of Fourier coefficients retained.

Since $m \ll d$, this is a considerable improvement on the previous estimator. However, some accuracy is lost through the transformation of the messages between the original and Fourier domains.

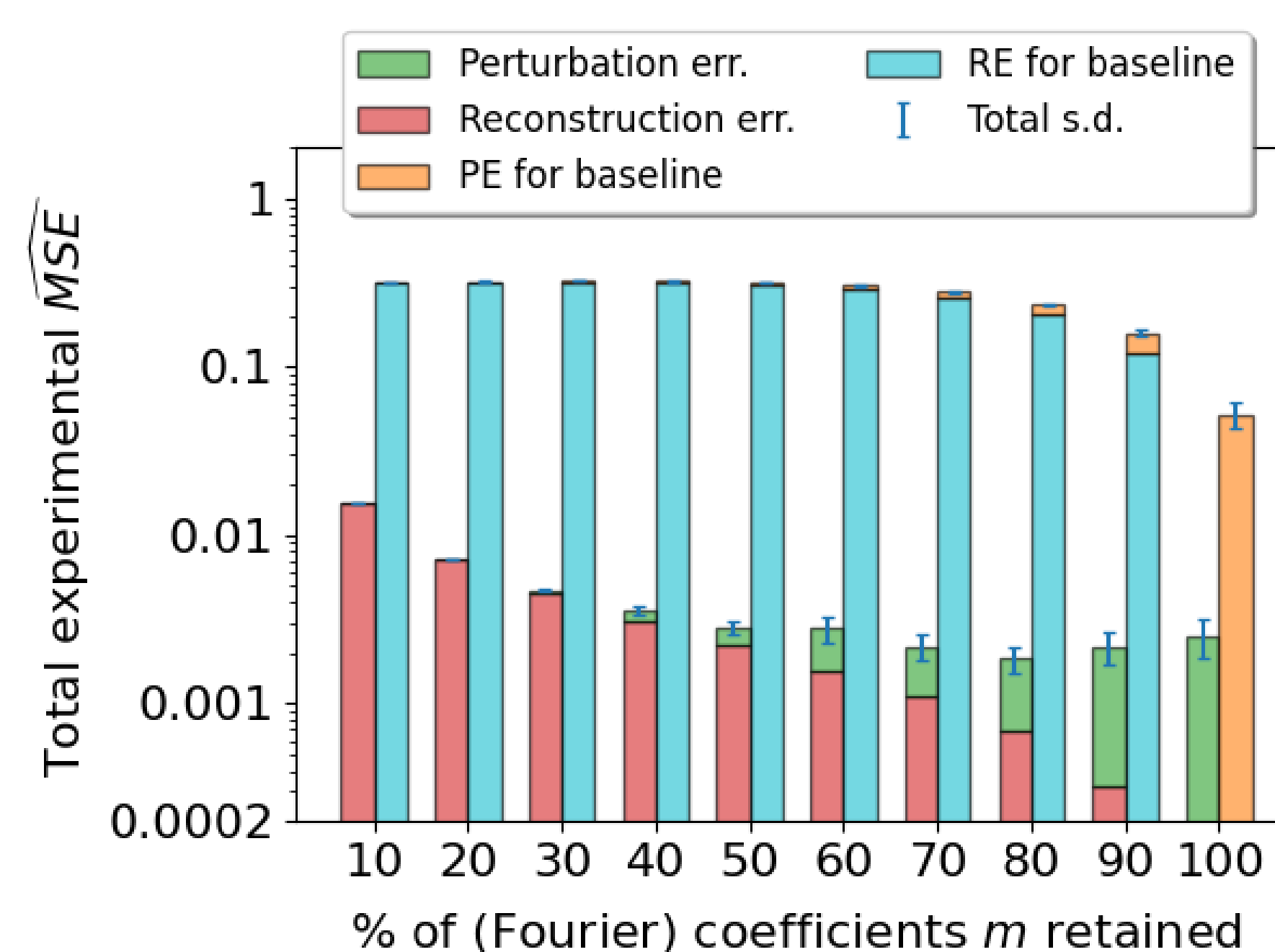


Figure 2: Relationship between experimental errors for a synthetic dataset created in Python.

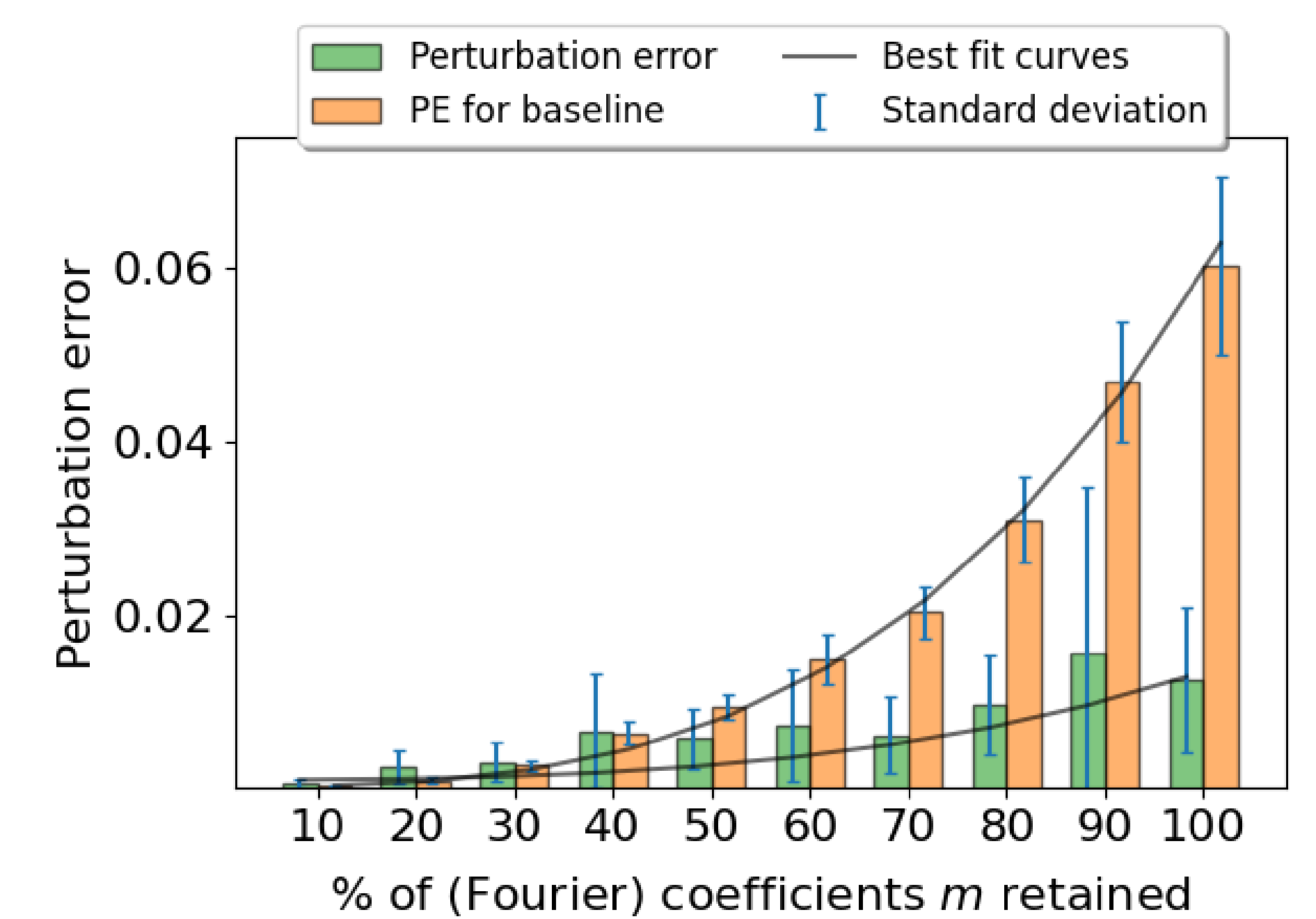


Figure 3: Perturbation error for the ECG Heartbeat Categorization Dataset, with best fit curve confirming $m^{8/3}$ dependency.

Experimental Evaluation

In the experimental section, we test our earlier hypotheses related to both protocols via a Python implementation, where the vectors are taken from an ECG Heartbeat Categorization Dataset. Initially, we analyse the effect of changing one key parameter at a time, while the others remain the same.

Figure 2 displays the effect of changing the number of Fourier coefficients m on the ratio between the perturbation and reconstruction errors. To illustrate this pattern more clearly, this graph has been plotted using a randomly generated synthetic dataset with a sinusoidal dependence on each coordinate.

Note that in Figure 2 the scale on the y-axis is logarithmic. Although the perturbation error initially looks much larger than its baseline counterpart, isolating these errors in Figure 3 clearly shows that the opposite is true for the ECG Heartbeat Categorization Dataset. The same conclusion can be made for the synthetic dataset.

Conclusion

The experiments above confirm that picking $t = 1$ and $k = 3$ serves to minimize the error. The lines of best fit, in Figure 3 for example, confirm the dependencies on the parameters m , d , ϵ and n .

We have seen via both theory and experiments that combining our new private summation protocol with a DFT reduces the MSE significantly, from a dependence on $d^{8/3}$ to $m^{8/3}$.

References

- [1] C. Dwork. Differential privacy. *Proc. 33rd Int. Colloq. Automata, Lang. Program. (ICALP)*, pages 1–12, 2006.
- [2] A. Cheu, A. Smith, J. Ullman, D. Zeber, and M. Zhilyaev. Distributed differential privacy via shuffling. *Proc. Annu. Int. Conf. Theory Appl. Cryptograph. Techn.*, pages 375–403, 2019.
- [3] B. Balle. The privacy blanket of the shuffle model. *Proc. Annu. Int. Cryptol. Conf.*, pages 638–667, 2019.
- [4] V. Rastogi and S. Nath. Differentially private aggregation of distributed time-series with transformation and encryption. In *Proc. ACM SIGMOD Int. Conf. Manage. Data*, pages 735–746. ACM, 2019.