

Wyner-Ziv Estimators for Distributed Mean Estimation with Side Information and Optimization

Prathamesh Mayekar
National University of Singapore
pratha22@nus.edu.sg

Ananda Theertha Suresh
Google Research
theertha@google.com

Shubham Jha
Indian Institute of Science
shubhamkj@iisc.ac.in

Himanshu Tyagi
Indian Institute of Science
htyagi@iisc.ac.in

Abstract

Communication efficient distributed mean estimation is an important primitive that arises in many distributed learning and optimization scenarios such as federated learning. Without any probabilistic assumptions on the underlying data, we study the problem of distributed mean estimation where the server has access to side information. We propose *Wyner-Ziv estimators*, which are communication and computationally efficient and near-optimal when an upper bound for the distance between the side information and the data is known. As a corollary, we also show that our algorithms provide efficient schemes for the classic Wyner-Ziv problem in information theory. In a different direction, when there is no knowledge assumed about the distance between side information and the data, we present an alternative Wyner-Ziv estimator that uses correlated sampling. This latter setting offers *universal recovery guarantees*, and perhaps will be of interest in practice when the number of users is large and keeping track of the distances between the data and the side information may not be possible.

With this mean estimator at our disposal, we revisit basic problems in decentralized optimization and compression where our Wyner-Ziv estimator yields algorithms with almost optimal performance. First, we consider the problem of communication constrained distributed optimization and provide an algorithm which attains the optimal convergence rate by exploiting the fact that the gradient estimates are close to each other. Specifically, the gradient compression scheme in our algorithm first uses half of the parties to form side information and then uses our Wyner-Ziv estimator to compress the remaining half of the gradient estimates.

Finally, we apply our Wyner-Ziv estimators to the classic Wyner-Ziv compression problem in information theory to get compression schemes that are computationally efficient and are almost optimal under much more relaxed assumptions than the standard probabilistic setting.

This work was supported by a grant from Robert Bosch Center for Cyber Physical Systems, Indian Institute of Science, and a grant on Security and Privacy for Smart Cities sponsored by National Security Council, India.

Parts of this paper appeared in the proceedings of International Conference on Artificial Intelligence and Statistics 2021 and the IEEE International Symposium on Information Theory 2022 ([41] and [40], respectively).

Contents

1	Introduction	4
1.1	Background	4
1.2	The model	4
1.3	Our contributions	6
1.4	Prior work	7
2	Preliminaries and the structure of our protocols	9
3	Distributed mean estimation with known Δ	9
3.1	Modulo Quantizer (MQ)	10
3.2	Rotated Modulo Quantizer (RMQ)	11
3.3	Subsampled RMQ: A Wyner-Ziv quantizer for \mathbb{R}^d	12
3.4	Lower bound	14
4	Distributed mean estimation for unknown Δ	15
4.1	The correlated sampling idea	15
4.2	Distance Adaptive Quantizer (DAQ)	15
4.3	Rotated Distance Adaptive Quantizer (RDAQ)	16
4.4	Subsampled RDAQ: A universal Wyner-Ziv quantizer for unit Euclidean ball	18
5	Application: Communication constrained distributed optimization	19
5.1	Lower bound	21
5.2	A general convergence bound	21
5.3	Baseline scheme: Parallel SGD	22
5.4	WZ-SGD: An almost optimal algorithm for distributed optimization	22
5.5	UWZ-SGD: A universal Wyner-Ziv algorithm for distributed optimization	24
6	The Gaussian Wyner-Ziv problem	26
7	The high-precision regime	27
7.1	RMQ in the high-precision regime.	27
7.2	Boosted RDAQ: RDAQ in the high-precision regime.	28
8	Numerical Experiments	30
9	Proofs	31
9.1	Proof of Lemma 2.1	31
9.2	Proof of Lemma 3.1	32
9.3	Proof of Lemma 3.2	33
9.4	Proof of Lemma 3.3	35
9.5	Proof of Theorem 3.5	36
9.6	Proof of Lemma 4.1	37
9.7	Proof of Lemma 4.2	38
9.8	Proof of Lemma 4.3	40
9.9	Proof of Theorem 5.1	41
9.10	Proof of Lemma 5.2	41

9.11 Proof of Theorem 5.3	42
9.12 Proof of Theorem 5.4	43
9.13 Proof of Theorem 5.5	49
9.14 Proof of Theorem 6.1	50
9.15 Proof of Lemma 7.2	51

1 Introduction

1.1 Background

Consider the problem of distributed mean estimation for n vectors $\{x_i\}_{i=1}^n$ in \mathbb{R}^d , where x_i is available to client i . Each client communicates to a server using a few bits to enable the server to compute the empirical mean

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i. \quad (1)$$

This estimation problem has become a crucial primitive for distributed optimization scenarios such as federated learning, where the data is distributed across multiple clients (see [11], [27], [30], [8], [48], [17], [25], [20], [51], [34], [49] [19], [10], [52], [58], [53], [60], [59], [38], [57], [4]). One of the main bottlenecks in such distributed scenarios is the significant communication cost incurred due to client communication at each iteration of the distributed algorithm. This has spurred a recent line of work which seeks to design quantizers to express x_i s using a low precision and, yet, enable the server to compute a high accuracy estimate of \bar{x} (see [55], [29], [14], [23], [43], [50], [7], [24], and the references therein).

Most of the recent works on distributed mean estimation focus on the setting where the server must estimate the sample mean based on the client vectors, and nothing else. However, in practice, the server may also have access to some side information. For example, consider the task of training a machine learning model based on remote client data as well as some publicly accessible data [9]. At each iteration, the server communicates its global model to the client, based on which the clients compute their updates (the gradient estimates based on their local data), compress them, and then send them to the server. The server may choose to compute its own update using the publicly available dataset to complement the updates from the client. In a related setting, the server can use the previously received gradients as side information for the next gradients expected from the clients. Alternatively, the server may ‘simulate’ side information from some client updates. It can then use this side information to form much more accurate estimates of other clients’ updates, leading to a faster distributed training algorithm. We discuss this application in detail in Section 5. Similarly, distributed mean estimation with side information can be used for variance reduction in other problems such as power iteration or parallel SGD (*cf.* [15]).

Motivated by these observations, for the distributed mean estimation problem described at the start of the section, we study the setting in which the server has access to the side information $\{y_i\}_{i=1}^n$ in \mathbb{R}^d , in addition to the communication from clients. Here, y_i can be viewed as server’s initial estimate (guess) of x_i . We emphasize that the side information y_i is available only to the sever and can, therefore, be used for estimating the mean at the server, but is not available to the clients while quantizing the updates $\{x_i\}_{i=1}^n$.

1.2 The model

Consider the input $\mathbf{x} := (x_1, \dots, x_n)$ and the side information $\mathbf{y} := (y_1, \dots, y_n)$. The clients use a communication protocol to send r bits each about their observed vector to the server. For the ease of implementation, we restrict to non-interactive protocols. Specifically, we allow *simultaneous message passing* (SMP) protocols $\pi = (\pi_1, \dots, \pi_n)$ where the communication $C_i = \pi_i(x_i, U) \in \{0, 1\}^r$

of client¹ i , $i \in [n]$, can only depend on its local observation x_i and public randomness U . Note that the clients are not aware of side information \mathbf{y} , which is available only to the server. In effect, the message C_i is obtained by *quantizing* x_i using an appropriately chosen randomized quantizer. Denoting the overall communication by $C^n := (C_1, C_2, \dots, C_n)$, the server uses the transcript (C^n, U) of the protocol and the side information \mathbf{y} to form the estimate of the sample mean² $\hat{x} = \hat{x}(C^n, U, \mathbf{y})$; see Figure 1 for a depiction of our setting. We call such a π an r -bit SMP protocol with input (\mathbf{x}, \mathbf{y}) and output \hat{x} .

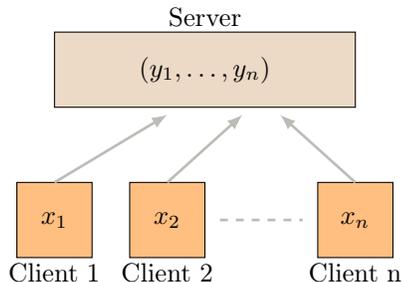


Figure 1: Problem setting of mean estimation with side information

We measure the performance of protocol π for inputs \mathbf{x} and \mathbf{y} and output \hat{x} using mean squared error (MSE) given by

$$\text{MSE}(\pi, \mathbf{x}, \mathbf{y}) := \mathbb{E} \left[\|\hat{x} - \bar{x}\|_2^2 \right],$$

where the expectation is over the public randomness U and \bar{x} is given in (1). We study the MSE of protocols for \mathbf{x} and \mathbf{y} such that the Euclidean distance between x_i and y_i is at most Δ_i , i.e.,

$$\|x_i - y_i\|_2 \leq \Delta_i, \quad \forall i \in [n]. \quad (2)$$

Denoting $\Delta := (\Delta_1, \dots, \Delta_n)$, we are interested in the performance of our protocols for the following two settings:

1. **The known Δ setting**, where Δ_i is known to client i and the server;
2. **The unknown Δ setting**, where Δ_i s are unknown to everyone.

In both these settings, we seek to find efficient r -bit quantizers for x_i that will allow accurate sample mean estimation. In the known Δ setting, the quantizers of different clients can be chosen using the knowledge of Δ ; in the unknown Δ setting, they must be fixed irrespective of Δ .

In another direction, we distinguish the *low-precision* setting of $r \leq d$ from the *high-precision* setting of $r > d$. The former is perhaps of more relevance for federated learning and high-dimensional distributed optimization, while the latter has received a lot of attention in the information theory literature on rate-distortion theory. Moreover, the distributed estimation problem is a lot more interesting in the low-precision setting. We, therefore, focus more on this regime while also providing extensions of our protocols to the high-precision regime.

As a benchmark, we recall the result for distributed mean estimation with no side-information from [55]. When all x_i s lie in the Euclidean ball of radius 1, [55] showed that the minmax MSE in

¹ $[n] := \{1, \dots, n\}$.

²While side information y_i is associated with client i , we do not enforce this association in our general formulation at this point.

the no side-information case is

$$\Theta\left(\frac{d}{nr}\right). \quad (3)$$

1.3 Our contributions

Drawing on ideas from distributed quantization problem in information theory (*cf.* [61]), specifically the Wyner-Ziv problem, we present *Wyner-Ziv estimators* for distributed mean estimation. In the known Δ setting, for a fixed Δ , and the low-precision setting of $r \leq d$, we propose an *r-bit SMP protocol* $\pi_{\mathbf{k}}^*$ which satisfies³

$$\text{MSE}(\pi_{\mathbf{k}}^*, \mathbf{x}, \mathbf{y}) = O\left(\sum_{i=1}^n \frac{\Delta_i^2}{n} \cdot \frac{d \log \log n}{nr}\right),$$

for all \mathbf{x} and \mathbf{y} satisfying (2). Thus, in the case where all x_i s lie in the Euclidean ball of radius 1, we improve upon the optimal estimator for distributed mean estimation (3) in the regime $\sum_{i=1}^n \frac{\Delta_i^2 \log \log n}{n} \leq 1$. Our estimator is motivated by the classic Wyner-Ziv problem, and hence, we refer to it as the *Wyner-Ziv estimator*. The details of the algorithm are given in Section 3.3.

Our protocol uses the same (randomized) r -bit quantizer for each client's data and simply uses the sample mean of the quantized vectors as the estimate for \bar{x} . Furthermore, the common quantizer used by the clients is efficient and has nearly linear time-complexity of $O(d \log d)$. Our proposed quantizer first applies a random rotation (proposed in [6]) to the input vectors x_i at client i and the side information vector y_i at the server. This ensures that the Δ_i upper bound on the ℓ_2 distance of x_i and y_i is converted to roughly a Δ_i/\sqrt{d} upper bound on the ℓ_∞ distance between x_i and y_i . This then enables us to use efficient one-dimensional quantizers for each coordinate of the x_i , which can now operate with the knowledge that the server knows a y_i with each coordinate within roughly Δ_i/\sqrt{d} of x_i 's coordinates.

Moreover, we show that this protocol $\pi_{\mathbf{k}}^*$ has optimal (worst-case) MSE up to an $O(\log \log n)$ factor. That is, we show that for any other r -bit SMP protocol π for $r \leq d$, we can find \mathbf{x} and \mathbf{y} satisfying (2) such that

$$\text{MSE}(\pi, \mathbf{x}, \mathbf{y}) = \Omega\left(\min_{i \in \{1, \dots, n\}} \Delta_i^2 \cdot \frac{d}{nr}\right).$$

In the unknown Δ setting, we propose a protocol $\pi_{\mathbf{u}}^*$ which adapts to the unknown distance Δ_i between x_i and y_i and, remarkably, provides MSE guarantees dependent on Δ . Specifically, for the low-precision setting of $r \leq d$, the protocol satisfies⁴

$$\text{MSE}(\pi_{\mathbf{u}}^*, \mathbf{x}, \mathbf{y}) = O\left(\sum_{i=1}^n \frac{\Delta_i}{n} \cdot \frac{d \log^* d}{nr}\right),$$

for all \mathbf{x} and \mathbf{y} in the unit Euclidean ball $\mathcal{B} := \{x \in \mathbb{R}^d : \|x\|_2 \leq 1\}$ and satisfying (2). Thus, we improve upon the optimal estimator for the no side information counterpart (3) in the regime $\sum_{i=1}^n \frac{\Delta_i \ln^* d}{n} \leq 1$. Once again, the quantizer employed by the protocol is efficient and has nearly

³We denote by $\log(\cdot)$ logarithm to the base 2 and by $\ln(\cdot)$ logarithm to the base e .

⁴We denote by $\ln^*(a)$ the minimum number of iterated logarithms to the base e that must be applied to a to make it less than 1.

linear time-complexity of $O(d \log d)$. At the heart of our proposed quantizer is the technique of correlated sampling from [21] which enables to derive a Δ dependent MSE bound.

Furthermore, both our quantizers can be extended to the high-precision regime of $r > d$. The quantizer for the known Δ setting directly extends by using r/d bits per dimension. The MSE of the SMP protocol using this quantizer for all the clients is only a factor of $\log n + r/d$ from the lower bound derived in [15] for the high-precision regime. The quantizer for the unknown Δ setting can be extended by sending the “type” of the communication vector, following an idea proposed in [42]. The MSE of the SMP protocol using this quantizer for all the clients falls as $2^{-r/d \ln^* d}$ as opposed to d/r that can be obtained using naive extensions of our quantizer.

As remarked at the outset, mean estimator is a basic primitive that can be used in problems related to decentralized optimization. Indeed, we apply our Wyner-Ziv estimator to a basic communication-constrained optimization problem and show that it leads to much faster algorithms for communication-constrained distributed optimization. Our first algorithm WZ-SGD significantly improves over the baseline Parallel SGD algorithm and is almost optimal for a large number of remote clients. We also propose a universal distributed optimization algorithm UWZ-SGD, where the remote clients can operate without the knowledge of the stochastic gradient’s variance. UWZ-SGD, too, improves the performance of the baseline Parallel SGD algorithm for large enough remote clients.

Finally, in a different direction, we revisit the classic Gaussian rate-distortion problem (*cf.* [46]) in information theory. In this problem, the encoder observing an Gaussian vector X wants to send it to a decoder observing a correlated Gaussian vector Y using r bits. Using the quantizer developed in the known Δ setting, we obtain an efficient scheme for this classic problem which requires a minuscule excess rate over the optimal asymptotic rate. Our scheme for this classic problem is interesting for two reasons: The first that it gives almost optimal result while using “covering” for each coordinate separately and hence is computationally efficient. All the existing schemes rely on high-dimensional covering constructed using structured codes and most of them are computationally inefficient. The second reason is that we do not require the distribution to be exactly Gaussian and subgaussianity suffices.

1.4 Prior work

The known Δ setting described above was first considered in [15]. The scheme of [15] relies on lattice quantizers with information theoretically optimal covering radius. Explicit lattices to be used and computationally efficient decoding is not provided.

In contrast, we provide explicit computationally efficient protocols for both low- and high-precision settings. Also, we establish lower bounds showing the optimality of our quantizer upto a multiplicative factor of $\log \log n$ in the low-precision regime of $r \leq d$. In comparison, the scheme of [15] is off by a factor of $\frac{d}{r}$ from this lower bound. Thus, when $r \ll d$, our scheme performs significantly better than that in [15]. We remark that the unknown Δ setting, which is perhaps more important in certain applications where estimating the distance of side information of each client is infeasible, has not been considered before.

In the classic information theoretic setting, related problems of quantization with side information at the decoder have been considered in rate-distortion theory starting with the seminal work of Wyner and Ziv [61]. Practical codes for settings where the observations are generated from known distributions have been constructed using channel codes; see, for instance, [31, 35, 37, 47, 62]. However, these codes are computationally too expensive for our setting, cannot be directly used

for our distribution-free setup, and are designed for the high-precision setting of $r > d$. We remark that the scheme proposed in [15] is similar to lattice schemes in [35, 37, 62].

The version of the distributed mean estimation problem with no side information at the server has been extensively studied. For any protocol in this setting operating with a precision constraint of $r \leq d$ bits per client, using a strong data processing inequality from [16], [55] shows a lower bound on MSE of $\Omega\left(\frac{d}{nr}\right)$, when all x_i s lie in the Euclidean ball of radius one. [55] propose a rotation based uniform quantization scheme which matches this lower bound up to a factor of $\log \log d$ for any precision constraint r . This upper bound is further improved by a random rotation based adaptive quantizer in [43] to a much tighter $\log \log^* d$ factor. For a precision constraint of $r = \Theta(d)$, the variable-length quantizers proposed in [55], [8], [48] as well as the fixed-length quantizers in [42], [19] are order-wise optimal.

A recent work on distributed mean estimation [26], which came after the conference version of our paper [41], proposed two different schemes for distributed mean estimation. The first scheme improves the performance of the standard Rand-k (*cf.* [53]) estimator when data across the clients are correlated. The second scheme uses previous gradient updates to improve the performance of the standard scheme. Using previous gradient updates can be seen as a special case of our setup when we use a historical gradient as side information. Interestingly, the second scheme in [26] uses the idea of centering the gradient estimate around the side information [26, Equation 12], which is similar to the decoding rule used in our second Wyner-Ziv estimate (14). A follow-up work of [41], [32] also proposed using correlation amongst clients to improve over standard sample mean estimators. Another recent work [54], which also came up after our conference version [41], proposed using correlated randomness for stochastic quantization across clients to improve the performance of the standard scheme.

[33] and an application considered in [54] are closest to the application of communication distributed optimization considered in Section 5. [33] builds on the distributed mean estimation schemes in [41] and proposes an algorithm for non-convex distributed optimization. However, unlike our proposed schemes, [33] suggests using historical gradients as side information, and its optimality is unclear.

[54] considers the same setting for communication-constrained distributed optimization as considered in this paper. The proposed scheme, too, is similar to UWZ – SGD, one of the schemes proposed in this paper. In more detail, both schemes leverage the fact that the stochastic estimates of the gradients across clients are close to each other to reduce the compression error. Moreover, they do this by using correlated randomness, and the compression can operate without knowing how close the stochastic gradients are across clients. However, there are crucial differences between the two schemes. At a high level, our scheme is designed for the low precision setting (where per client precision is less than the dimension) and only uses a fixed length code, the scheme in [54] is designed for the high precision setting and uses a variable length code in this setting.

Our results for the low-precision regime in known Δ setting are provided in Section 3 and in the unknown Δ setting are provided in Section 4. In Section 7, we extend our results to the high-precision regime. In Section 5, we derive new algorithms for communication-constrained distributed optimization using our distributed mean estimation protocols. In Section 6, we provide an application of the quantizer developed for the known-setting to the Gaussian Wyner-Ziv problem. Finally, we close with all the proofs in Section 9. Before presenting these results, we review some preliminaries in the next section.

2 Preliminaries and the structure of our protocols

While our lower bound for the known Δ setting holds for an arbitrary SMP protocol, both the protocols we propose in this paper, for the known Δ and the unknown Δ settings, have a common structure. We use r -bit quantizers to form estimates of x_i s at the server and then compute the sample mean of the estimates of x_i s. To describe our protocols and facilitate our analysis, we begin by concretely defining the distributed quantizers needed for this problem. Further, we present a simple result relating the performance of the resulting protocol to the parameters of the quantizer.

An r -bit quantizer Q for input vectors in $\mathcal{X} \subset \mathbb{R}^d$ and side information $\mathcal{Y} \subset \mathbb{R}^d$ consists of randomized mappings⁵ (Q^e, Q^d) with the encoder mapping $Q^e : \mathcal{X} \rightarrow \{0, 1\}^r$ used by the client to quantize and the decoder mapping $Q^d : \{0, 1\}^r \times \mathcal{Y} \rightarrow \mathcal{X}$ used by the server to aggregate quantized vectors. The overall quantizer Q is given by the composition mapping $Q(x, y) = Q^d(Q^e(x), y)$.

In our protocols, for input \mathbf{x} and side information \mathbf{y} , client i uses the encoder Q_i^e for the r -bit quantizer Q_i to send $Q_i^e(x_i)$. The server uses $Q_i^e(x_i)$ and y_i to form the estimate $\hat{x}_i = Q_i(x_i, y_i)$ of x_i . We assume that the randomness used in quantizers Q_i for different i is independent, whereby \hat{x}_i are independent of each other for different i . Then server finally forms the estimate of the sample mean as

$$\hat{\bar{x}} := \frac{1}{n} \sum_{i=1}^n \hat{x}_i. \quad (4)$$

For any quantizer Q , the following two quantities will determine its performance when used in our distributed mean estimation protocol:

$$\begin{aligned} \alpha(Q; \Delta) &:= \sup_{x \in \mathcal{X}, y \in \mathcal{Y}: \|x-y\|_2 \leq \Delta} \mathbb{E} [\|Q(x, y) - x\|_2^2], \\ \beta(Q; \Delta) &:= \sup_{x \in \mathcal{X}, y \in \mathcal{Y}: \|x-y\|_2 \leq \Delta} \|\mathbb{E} [Q(x, y) - x]\|_2^2, \end{aligned}$$

where the expectation is over the randomization of the quantizer. Note that $\alpha(Q; \Delta)$ can be interpreted as the worst-case MSE and $\beta(Q, \Delta)$ the worst-case bias over $x \in \mathcal{X}$ and $y \in \mathcal{Y}$ such that $\|x - y\|_2 \leq \Delta$.

The result below will be very handy for our analysis.

Lemma 2.1. *For $\mathbf{x} \in \mathcal{X}^n$ and $\mathbf{y} \in \mathcal{Y}^n$ satisfying (2) and r -bit quantizers Q_i , $i \in [n]$, using independent randomness for different $i \in [n]$, the estimate $\hat{\bar{x}}$ in (4) and the sample mean \bar{x} in (1) satisfy*

$$\mathbb{E} [\|\hat{\bar{x}} - \bar{x}\|_2^2] \leq \sum_{i=1}^n \frac{\alpha(Q_i; \Delta_i)}{n^2} + \sum_{i=1}^n \frac{\beta(Q_i; \Delta_i)}{n}.$$

3 Distributed mean estimation with known Δ

In this section, we present our Wyner-Ziv estimator for the known Δ setting. As described in Section 2, we use the the same (randomized) quantizer across all the clients and form the estimate of sample mean as in (4). We only need to define the common quantizer used by all the clients, which we do in Section 3.3. In Sections 3.1 and 3.2, we provide the basic building blocks of our final

⁵We can use public randomness U for randomizing.

quantizer. Further, in Section 3.4, we derive a lower bound for the worst-case MSE that establishes the near-optimality of our protocol. Throughout we restrict to the low-precision setting of $r \leq d$.

3.1 Modulo Quantizer (MQ)

The first subroutine used by our larger quantizer is the *Modulo Quantizer* (MQ). MQ is a one dimensional distributed quantizer that can be applied to the input $x \in \mathbb{R}$ with side information $y \in \mathbb{R}$. We give an input parameter Δ' to MQ where $|x - y| \leq \Delta'$. In addition to Δ' , MQ also has the resolution parameter k and the lattice parameter ε as inputs.

For an appropriate ε to be specified later, we consider the lattice $\mathbb{Z}_\varepsilon = \{\varepsilon z : z \in \mathbb{Z}\}$. For a given input x , the encoder Q_M^e finds the closest points in \mathbb{Z}_ε larger and smaller than x . Then, one of these points is sampled randomly to get an unbiased estimate of x . The sampled point will be of the form $\tilde{z}\varepsilon$, where \tilde{z} is in \mathbb{Z} . We note that the chosen point \tilde{z} satisfies

$$\begin{aligned} \varepsilon \mathbb{E}[\tilde{z}] &= x \text{ and} \\ |x - \varepsilon \tilde{z}| &< \varepsilon, \quad \text{almost surely.} \end{aligned} \tag{5}$$

The encoder sends $w = \tilde{z} \bmod k$ to the decoder, which requires $\log k$ bits.

Upon receiving this w , the decoder Q_M^d looks at the set $\mathbb{Z}_{w,\varepsilon} = \{(zk + w) \cdot \varepsilon : z \in \mathbb{Z}\}$ and decodes the point closest to y , which we denote by $Q_M(x, y)$. Note that declaring y will already give a MSE of less than Δ . A useful property of this decoder is that its output is always within a bounded distance from y ; namely, since in Step 1 of Alg. 3 we look for the closest point to y in the lattice $\mathbb{Z}_{w,\varepsilon} := \{(zk + w) \cdot \varepsilon : z \in \mathbb{Z}\}$, the output $Q_M(x, y)$ satisfies

$$|Q_M(x, y) - y| \leq k\varepsilon, \quad \text{almost surely.} \tag{6}$$

We summarize MQ in Alg. 2 and 3.

Require: Input $x \in \mathbb{R}$, Parameters k, Δ' , and ε

- 1: Compute $z_u = \lceil x/\varepsilon \rceil, z_l = \lfloor x/\varepsilon \rfloor$
- 2: Generate $\tilde{z} = \begin{cases} z_u, & \text{w.p. } x/\varepsilon - z_l \\ z_l, & \text{w.p. } z_u - x/\varepsilon \end{cases}$
- 3: **Output:** $Q_M^e(x) = \tilde{z} \bmod k$

Algorithm 2: Encoder $Q_M^e(x)$ of MQ

Require: Input $w \in \{0, \dots, k-1\}, y \in \mathbb{R}$

- 1: Compute $\hat{z} = \arg \min\{|(zk + w) \cdot \varepsilon - y| : z \in \mathbb{Z}\}$
- 2: **Output:** $Q_M^d(w, y) = (\hat{z}k + w)\varepsilon$

Algorithm 3: Decoder $Q_M^d(w, y)$ of MQ

The result below provides performance guarantees for Q_M . The key observation is that the output $Q_M(x, y)$ of the quantizer equals $\tilde{z}\varepsilon$ with \tilde{z} found at the encoder, if ε is set appropriately.

Lemma 3.1. Consider the Modulo Quantizer Q_M described in Alg. 2 and 3 with parameter ε set to satisfy

$$k\varepsilon \geq 2(\varepsilon + \Delta'). \quad (7)$$

Then, for every x, y in \mathbb{R} such that $|x - y| \leq \Delta'$, the output $Q_M(x, y)$ of MQ satisfies

$$\begin{aligned} \mathbb{E}[Q_M(x, y)] &= x \quad \text{and} \\ |Q_M(x, y) - x| &\leq \varepsilon, \quad \text{almost surely.} \end{aligned}$$

In particular, we can set $\varepsilon = 2\Delta'/(k - 2)$, to get $|Q_M(x, y) - x| \leq 2\Delta'/(k - 2)$. Furthermore, the output of Q_M can be described in $\log k$ bits.

We close with a remark that the modulo operation used in our scheme is the simplest and easily implementable version of classic coset codes obtained using nested lattices used in distributed quantization (cf. [18, 36, 62]) and was used in [15] as well.

3.2 Rotated Modulo Quantizer (RMQ)

We now describe *Rotated Modulo Quantizer (RMQ)*. RMQ and the subsequent quantizers in this section will be used to quantize input vector x in \mathbb{R}^d with side information y in \mathbb{R}^d , where $\|x - y\|_2 \leq \Delta$. RMQ first preprocesses the input x and side information y by randomly rotating them and then simply applies MQ for each coordinate. For rotation, we multiply both x and y with a matrix R given by

$$R = \frac{1}{\sqrt{d}} \cdot HD, \quad (8)$$

where H is the $d \times d$ Walsh-Hadamard Matrix (see [22])⁶ and D is a diagonal matrix with each diagonal entry generated uniformly from $\{-1, +1\}$. Note that we use public randomness⁷ to generate the same D at both the encoder and the decoder. We formally describe the quantizer in⁸ Alg. 4 and 5.

Remark 1. We remark that the vector $R(x - y)$ has zero mean subgaussian coordinates with a variance factor of Δ^2/d . This implies that for all coordinates i in $[d]$, we have

$$P(|R(x - y)(i)| \geq \Delta') \leq 2e^{-\frac{\Delta'^2 d}{2\Delta^2}}$$

(see, for instance, [12, Theorem 2.8]). This observation allows us to use $\Delta' \approx \Delta/\sqrt{d}$ for MQ applied to each coordinate.

Lemma 3.2. Fix $\Delta \geq 0$. Let $Q_{M,R}$ be RMQ described in Alg. 4 and 5. Then, for⁹ $k \geq 4$, $\delta \in (0, \Delta)$, $\Delta' = \sqrt{6(\Delta^2/d) \ln(\Delta/\delta)}$ and the parameter ε of MQ set to $\varepsilon = 2\Delta'/(k - 2)$, we have for $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$ that

$$\alpha(Q_{M,R}; \Delta) \leq \frac{24 \Delta^2}{(k - 2)^2} \ln \frac{\Delta}{\delta} + 154 \delta^2 \quad \text{and}$$

⁶We assume that d is a power of 2. If it isn't, we can pad the vector by zeros to make it a power of 2; even in the worst-case, this only doubles the required bits.

⁷In practice, this can be implemented by using the same seed for pseudo-random number generator at encoder and decoder.

⁸We denote by (e_1, \dots, e_d) the standard basis of \mathbb{R}^d .

⁹In the proof, we provide a general bound which holds for all k .

Require: Input $x \in \mathbb{R}^d$, Parameters k and Δ'

- 1: Sample R as in (8) using public randomness
- 2: $x' = Rx$
- 3: **Output:** $Q_{M,R}^e(x) = [Q_M^e(x'(1)), \dots, Q_M^e(x'(d))]^T$ using parameters k, ε , and Δ' for Q_M^e of Alg. 2

Algorithm 4: Encoder $Q_{M,R}^e(x)$ of RMQ

Require: Input $w \in \{0, \dots, k-1\}^d, y \in \mathbb{R}^d$,
Parameters k and Δ'

- 1: Get R from public randomness.
- 2: $y' = Ry$
- 3: **Output:** $Q_{M,R}^d(w, y) = R^{-1} \sum_{i \in [d]} Q_M^d(w(i), y'(i))e_i$ using parameters k, ε , and Δ' for Q_M^d of Alg. 3,

Algorithm 5: Decoder $Q_{M,R}^d(w, y)$ of RMQ

$$\beta(Q_{M,R}; \Delta) \leq 154 \delta^2.$$

Furthermore, the output of quantizer $Q_{M,R}$ can be described in $d \log k$ bits.

Remark 2. The choice of Δ' in the first statement of the Lemma 3.2 is based on Remark 1. We note that δ is a parameter to control the bias incurred by our quantizer. By setting $\Delta' = \Delta$ we can get an unbiased quantizer, but it only recovers the performance obtained by simply using MQ for each coordinate, an algorithm considered in [15] as well.

3.3 Subsampled RMQ: A Wyner-Ziv quantizer for \mathbb{R}^d

Our final quantizer is a modification of RMQ of previous section where we make the precision less than r bits by randomly sampling a subset of coordinates. Specifically, note that $Q_{M,R}^e(x)$ sends d binary strings of $\log k$ bits each. We reduce the resolution by sending only a random subset S of these strings. This subset is sampled using shared randomness and is available to the decoder, too. Note that $Q_{M,R}^d$ applies Q_M^d to these strings separately; now, we use Q_M^d to decode the entries in S alone. We describe the overall quantizer in Alg. 6 and 7.

Require: Input $x \in \mathbb{R}$, Parameters k, Δ' , and μ

- 1: Sample $S \subset [d]$ u.a.r. from all subsets of $[d]$ of cardinality μd and sample R as in (8) using public randomness
- 2: **Output:** $Q_{wz}^e(x) = \{Q_M^e(Rx(i)) : i \in S\}$ using parameters k, ε , and Δ' for Q_M^e of Alg. 2

Algorithm 6: Encoder $Q_{wz}^e(x)$ of subsampled RMQ

Require: Input $w \in \{0, \dots, k-1\}^{\mu d}$, $y \in \mathbb{R}$

- 1: Get S and R from public randomness
- 2: Compute $\tilde{x} = (Q_M^d(w(i), Ry(i)), i \in S)$ using parameters k , ε , and Δ' for Q_M^d of Alg. 3
- 3: $\hat{x}_R = \frac{1}{\mu} \sum_{i \in S} (\tilde{x}(i) - Ry(i)) e_i + Ry$
- 4: **Output:** $Q_{wz}^d(w, y) = R^{-1} \hat{x}_R$

Algorithm 7: Decoder $Q_{wz}^d(w, y)$ of subsampled RMQ

Remark 3. We remark that, typically, when implementing random sampling, we set the unsampled components to 0. However, to get Δ dependent bounds on MSE, we set the unsampled coordinates to the corresponding coordinate of side information and center our estimate appropriately to only have small bias.

The result below relates the performance of our final quantizer Q_{wz} to that of $Q_{M,R}$, which was already analysed in the previous section.

Lemma 3.3. *Fix $\Delta > 0$. Let Q_{wz} and $Q_{M,R}$ be the quantizers described in Alg. 6 and 7 and Alg. 4 and 5, respectively. Then, for $\mu d \in [d]$, we have for $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$ that*

$$\alpha(Q_{wz}; \Delta) \leq \frac{2\alpha(Q_{M,R}; \Delta)}{\mu} + \frac{2\Delta^2}{\mu} \quad \text{and}$$

$$\beta(Q_{wz}; \Delta) = \beta(Q_{M,R}; \Delta).$$

Furthermore, the output of quantizer Q_{wz} can be described in $\mu d \log k$ bits.

We are now equipped to prove our first main result. Our protocol π_k^* uses Q_{wz} for each client as described in Section 2 and forms the estimate \hat{x} as in (4). We set the parameters needed for Q_{wz} in Alg. 6 and 7 as follows: For client i , we set the parameters of MQ as

$$\delta = \frac{\Delta_i}{\sqrt{n}}, \quad \log k = \left\lceil \log(2 + \sqrt{12 \ln n}) \right\rceil, \quad \Delta' = \sqrt{6(\Delta_i^2/d) \ln(\Delta_i/\delta)}, \quad \varepsilon = 2\Delta'/(k-2), \quad (9)$$

and set the parameter μ as

$$\mu d = \left\lfloor \frac{r}{\log k} \right\rfloor. \quad (10)$$

We characterize the resulting error performance in the next result.

Theorem 3.4. *For a $n \geq 2$, a fixed $\mathbf{\Delta} = (\Delta_1, \dots, \Delta_n)$, and $d \geq r \geq 2 \left\lceil \log(2 + \sqrt{12 \ln n}) \right\rceil$, the protocol π_k^* with parameters as set in (9) and (10) is an r -bit protocol which satisfies*

$$\text{MSE}(\pi_k^*, \mathbf{x}, \mathbf{y}) \leq (79 \lceil \log(2 + \sqrt{12 \ln n}) \rceil + 26) \left(\sum_{i=1}^n \frac{\Delta_i^2}{n} \cdot \frac{d}{nr} \right),$$

for all \mathbf{x}, \mathbf{y} satisfying (2).

Proof. Denoting by Q_i the quantizer Q_{WZ} with parameters set for user i , by Lemmas 2.1 and 3.3, we get

$$\begin{aligned} \mathbb{E} [\|\hat{x} - \bar{x}\|_2^2] &\leq \sum_{i=1}^n \frac{\alpha(Q_i; \Delta_i)}{n^2} + \sum_{i=1}^n \frac{\beta(Q_i; \Delta_i)}{n} \\ &\leq \frac{1}{\mu n^2} \sum_{i=1}^n (\alpha(Q_{\text{M},R,i}; \Delta_i) + \Delta_i^2) + \sum_{i=1}^n \frac{\beta(Q_{\text{M},R,i}; \Delta_i)}{n}, \end{aligned}$$

where $Q_{\text{M},R,i}$ denotes RMQ with parameters set for user i . Further, since $k \geq 4$ holds when $n \geq 2$ for our choice of parameters, by using Lemma 3.2 and substituting $\delta^2 = \Delta_i^2/n$, we get

$$\begin{aligned} \alpha(Q_{\text{M},R,i}; \Delta_i) &\leq \frac{12\Delta_i^2 \ln n}{(k-2)^2} + \frac{154\Delta_i^2}{n}, \\ \beta(Q_{\text{M},R,i}; \Delta_i) &\leq \frac{154\Delta_i^2}{n}, \end{aligned}$$

which with the previous bound gives

$$\begin{aligned} \mathbb{E} [\|\hat{x} - \bar{x}\|_2^2] &\leq \frac{1}{\mu d} \left(\frac{12 \ln n}{(k-2)^2} + \frac{154}{n} + 1 + 154\mu \right) \sum_{i=1}^n \frac{d\Delta_i^2}{n^2} \\ &\leq \frac{79 \lceil \log(2 + \sqrt{12 \ln n}) \rceil + 26}{r} \sum_{i=1}^n \frac{d\Delta_i^2}{n^2}, \end{aligned}$$

where in the final bound we used our choice of k , the assumption that $n \geq 2$ (which implies that $d \geq r \geq 6$), and the fact that $\lceil r/\log k \rceil \geq r/2$ if $r \geq 2 \log k$. \square

Remark 4. We note that by using MQ for each coordinate without rotating (or even with rotation using R as above) and with $\Delta' = \Delta_i$ yields MSE less than

$$O \left(\sum_{i=1}^n \frac{\Delta_i^2}{n} \cdot \frac{d \log d}{nr} \right),$$

for $r \leq d$. Thus, our approach above allows us to remove the $\log d$ factor at the cost of a (milder for large d) $\log \log n$ factor.

Thus, as can be seen from the lower bound presented in Theorem 3.5 below, our Wyner-Ziv estimator π_k^* is nearly optimal. Finally, Q_{WZ} can be efficiently implemented as both the encoding and decoding procedures have nearly-linear time complexity¹⁰ of $O(d \log d)$.

3.4 Lower bound

We now prove a lower bound on the MSE incurred by any SMP protocol using r bits per client. The proof relies on the strong data processing inequality in [16] and is similar in structure to the lower bound for distributed mean estimation without side-information in [55].

¹⁰The most expensive operation at both the encoder and decoder of this estimator is the Hadamard matrix multiplication operation, which requires $d \log d$ real operations.

Theorem 3.5. Fix $\Delta = (\Delta_1, \dots, \Delta_n)$. There exists a universal constant $c < 1$ such that for any r -bit SMP protocol π , with $r \leq cd$, there exists input $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{2d}$ satisfying (2) and such that

$$\text{MSE}(\pi, \mathbf{x}, \mathbf{y}) \geq c \min_{i \in [d]} \Delta_i^2 \cdot \frac{d}{nr}.$$

4 Distributed mean estimation for unknown Δ

Finally, we present our Wyner-Ziv estimator for the unknown Δ setting. We first, in Section 4.1, describe the idea of correlated sampling from [21], which will serve as an essential building block for all our quantizers in this section. We then build towards our final quantizer, described in 4.4, by first describing its simpler versions in Section 4.2 and 4.3. Once again, we restrict to the low-precision setting of $r \leq d$.

4.1 The correlated sampling idea

Suppose we have two numbers x and y lying in $[0, 1]$. A 1-bit unbiased estimator for x is the random variable $\mathbb{1}_{\{U \leq x\}}$, where U is a uniform random variable in $[0, 1]$. The variance of such an estimator is $x - x^2$. We consider a variant of this estimator given by:

$$\hat{X} = \mathbb{1}_{\{U \leq x\}} - \mathbb{1}_{\{U \leq y\}} + y, \tag{11}$$

where, like before, U is a uniform random variable in $[0, 1]$. Such an estimator still uses only 1-bit of information related to x . It is easy to check that this estimator unbiased estimator of x , namely $\mathbb{E}[\hat{X}] = x$. The variance of this estimator is given by

$$\text{var}(\hat{X}) = \mathbb{E}[(\hat{X} - x)^2] = |x - y| - (x - y)^2,$$

which is lower than that of the former quantizer when x is close to y . We build-on this basic primitive to obtain a quantizer with MSE bounded above by a Δ -dependent expression, without requiring the knowledge of Δ .

4.2 Distance Adaptive Quantizer (DAQ)

DAQ and subsequent quantizers in this Section will be described for input x and side information y lying in \mathbb{R}^d . The first component of our quantizer, DAQ, which uses (11) and incorporates the correlated sampling idea discussed earlier. Both the encoder and the decoder of DAQ use the same d uniform random variables $\{U(i)\}_{i=1}^d$ between $[-1, 1]$, which are generated using public randomness. At the encoder, each coordinate of vector x is encoded to the bit $\mathbb{1}_{\{U(i) \leq x(i)\}}$. At the decoder, using the bits received from the encoder, side information y , and the public randomness $\{U(i)\}_{i=1}^d$, we first compute bits $\mathbb{1}_{\{U(i) \leq y(i)\}}$ for each $i \in [d]$. Then, the estimate of x is formed as follows:

$$Q_D(x, y) = \sum_{i=1}^d (\mathbb{1}_{\{U(i) \leq x(i)\}} - \mathbb{1}_{\{U(i) \leq y(i)\}}) e_i + y.$$

We formally describe the quantizer in Alg. 8 and 9. The next result characterizes the performance for DAQ.

Require: Input $x \in \mathbb{R}^d$

- 1: Sample $U(i) \sim \text{Unif}[-1, 1], \forall i \in [d]$
- 2: $\tilde{x} = \sum_{i=1}^d \mathbb{1}_{\{U(i) \leq x(i)\}} \cdot e_i$
- 3: **Output:** $Q_{\mathbb{D}}^{\circ}(x) = \tilde{x}$, where \tilde{x} is viewed as binary vector of length d

Algorithm 8: Encoder $Q_{\mathbb{D}}^{\circ}(x)$ of DAQ

Require: Input $w \in \{0, 1\}^d, y \in \mathbb{R}^d$,

- 1: Get $U(i), \forall i \in [d]$, using public randomness
- 2: Set $\tilde{y} = \sum_{i=1}^d \mathbb{1}_{\{U(i) \leq y(i)\}} \cdot e_i$
- 3: **Output:** $Q_{\mathbb{D}}^{\Delta}(w, y) = 2(w - \tilde{y}) + y$, where w is viewed as a vector in \mathbb{R}^d

Algorithm 9: Decoder $Q_{\mathbb{D}}^{\Delta}(w, y)$ of DAQ

Lemma 4.1. *Let $Q_{\mathbb{D}}$ denote DAQ described in Algorithms 8 and 9. Then, for $\mathcal{X} = \mathcal{Y} = \mathcal{B}$ and every $\Delta > 0$, we have*

$$\alpha(Q_{\mathbb{D}}; \Delta) \leq 2\Delta\sqrt{d} \quad \text{and} \quad \beta(Q_{\mathbb{D}}; \Delta) = 0.$$

Furthermore, the output of quantizer $Q_{\mathbb{D}}$ can be described in d bits.

4.3 Rotated Distance Adaptive Quantizer (RDAQ)

Next, we proceed as for the known Δ setting and add a preprocessing step of rotating x and y using random matrix R of (8), which is sampled using shared randomness. We remark that here random rotation is used to exploit the subgaussianity of the rotated x and y , whereas in RMQ of previous section it was used to exploit the subgaussianity of $x - y$. After this rotation step, we proceed with a quantizer similar to DAQ, but we quantize each coordinate at multiple ‘‘scales.’’ We describe this step in detail below.

Using multiple scales. In DAQ, we considered each coordinate x to be anywhere between $[-1, 1]$ and used one uniform random variable for each coordinate. Now, we will use h independent uniform random variables for each coordinate, each corresponding to a different scale $[-M_j, M_j]$, $j \in \{0, 1, 2, \dots, h-1\}$. For convenience, we abbreviate $[h]_0 := \{0, 1, 2, \dots, h-1\}$.

Specifically, let $U(i, j)$ be distributed uniformly over $[-M_j, M_j]$, independently for different $i \in [d]$ and different $j \in [h]_0$. The values M_j s correspond to different scales and are set, along with h , as follows: For all $j \in [h]_0$,

$$M_j^2 := \frac{6}{d} \cdot e^{*j}, \quad \log h := \lceil \log(1 + \ln^*(d/6)) \rceil, \quad (12)$$

where e^{*j} denotes the j th iteration of e given by $e^{*0} := 1$, $e^{*1} := e$, $e^{*j} := e^{e^{*(j-1)}}$. All the dh uniform random variables are generated using public randomness and are available to both the encoder and the decoder.

The intervals $[-M_j, M_j]$ are designed to minimize the MSE of our quantizer by tuning its ‘‘resolution’’ to the ‘‘scale’’ of the input, and while still ensuring unbiased estimates. This idea of

using multiple intervals $[-M_j, M_j]$ for quantizing the randomly rotated vector is from [43], where it was used for the case with no side information.

Multiscale DAQ. After rotation, we proceed as in DAQ, except that we use different scale M_j for different coordinates. Ideally, for the i th coordinate, we would like to use $M_{z^*(i)}$, where $z^*(i)$ is the smallest index such that both $Rx(i)$ and $Ry(i)$ lie in $[-M_{z^*(i)}, M_{z^*(i)}]$. However, since y is not available to the encoder, we simply resort to sending the smallest value $z(i)$ which is the smallest index such that $Rx(i) \in [-M_{z(i)}, M_{z(i)}]$ and apply the encoder of DAQ h times to compress x at all scales, *i.e.*, we send h bits $(\mathbb{1}_{\{U(i,j) \leq Rx(i)\}}, j \in [h]_0)$.

Thus, the overall number of bits used by RDAQ's encoder is $d \cdot (h + \lceil \log h \rceil)$. At RDAQ's decoder, using $z(i)$, we compute the smallest index $z^*(i)$ containing both $Rx(i)$ and $Ry(i)$. In effect, the decoder emulates the decoder for DAQ applied to Ry , but for scale $M_{z^*(i)}$. The encoding and decoding algorithm of RDAQ are described in Alg. 10 and 11, respectively.

Require: Input $x \in \mathcal{B}$

- 1: Sample $U(i, j) \sim \text{Unif}[-M_j, M_j]$, $i \in [d], j \in [h]_0$, and sample R as in(8) using public randomness.
- 2: $x_R = Rx$
- 3: **for** $i \in [d]$ **do**
 $z(i) = \min\{j \in [h]_0 : |x_R(i)| \leq M_j\}$
- 4: **for** $j \in [h]_0$ **do**
 $\tilde{x}_j = \sum_{i=1}^d \mathbb{1}_{\{U(i,j) \leq x_R(i)\}} e_i$
- 5: **Output:** $Q_{D,R}^e(x) = ([\tilde{x}_0, \dots, \tilde{x}_{h-1}], z)$, where we view \tilde{x}_j s as binary vectors

Algorithm 10: Encoder $Q_{D,R}^e(x)$ at for RDAQ

Require: Input $(w, z) \in \{0, 1\}^{d \times h} \times [h]_0^d$ and $y \in \mathcal{B}$

- 1: Get $U(i, j)$, $i \in [d], j \in [h]_0$, and R using public randomness.
- 2: $y_R = Ry$
- 3: **for** $i \in [d]$ **do**
 $z'(i) = \min\{j \in \{[h]_0\} : |y_R(i)| \leq M_j\}$
 $z^*(i) = \max\{z(i), z'(i)\}$
- 4: $w' = \sum_{i=1}^d 2M_{z^*(i)} (w(i, z^*(i)) - \mathbb{1}_{\{U(i, z^*(i)) \leq y_R\}})$
- 5: $\hat{x}_R = w' + Ry$
- 6: **Output:** $Q_{D,R}^d(w, y) = R^{-1} \hat{x}_R$.

Algorithm 11: Decoder $Q_{D,R}^d(x)$ for RDAQ

Then, the quantized output $Q_{D,R}$ corresponding to input vector x and side-information y is

$$Q_{D,R}(x, y) = R^{-1} \left[\sum_{i=1}^d 2M_{z^*(i)} (\mathbb{1}_{\{U(i, z^*(i)) \leq Rx(i)\}} - \mathbb{1}_{\{U(i, z^*(i)) \leq Ry(i)\}}) + Ry \right].$$

We remark that since rotated coordinates $Rx(i)$ and $Ry(i)$ have subgaussian tails, with very high probability $M_{z^*(i)}$ will be much less than 1, which helps in reducing the overall MSE significantly. The performance of the algorithm is characterized below.

Lemma 4.2. *Let $Q_{D,R}$ be RDAQ described in Alg. 10 and 11. Then, for $\mathcal{X} = \mathcal{Y} = \mathcal{B}$ and every $\Delta > 0$, we have*

$$\alpha(Q_{D,R}; \Delta) \leq 16\sqrt{3}\Delta \quad \text{and} \quad \beta(Q_{D,R}; \Delta) = 0.$$

Furthermore, the output of quantizer Q can be described in $d(h + \log h)$ bits.

4.4 Subsampled RDAQ: A universal Wyner-Ziv quantizer for unit Euclidean ball

Finally, we bring down the precision of RDAQ to r , as before for the known Δ setting, by retaining the output of RDAQ for only coordinates $i \in S$, where S is generated uniformly at random from all subsets of $[d]$ of cardinality μd using public randomness. Specifically, we execute Alg. 10 and 11 with S replacing $[d]$ and multiplying w' in Step 4 of Alg. 11 by normalization factor of $d/|S|$. The output of the resulting encoder is given by

$$Q_{\text{wz},u}^e(x) = \{Q_{D,R}^e(x)(i) : i \in S\}, \quad (13)$$

where $Q_{D,R}^e(x)(i)$ represents the encoded bits $([\tilde{x}_0(i), \dots, \tilde{x}_{h-1}(i)], z(i))$ for the i th coordinate using RDAQ, and the output of the resulting decoder is given by

$$Q_{\text{wz},u}(x, y) = R^{-1} \left[\frac{1}{\mu} \sum_{i \in S} 2M_{z^*(i)} \left(\mathbb{1}_{\{U(i, z^*(i)) \leq Rx(i)\}} - \mathbb{1}_{\{U(i, z^*(i)) \leq Ry(i)\}} \right) + Ry \right]. \quad (14)$$

Lemma 4.3. *Let $Q_{\text{wz},u}$ be the quantizers described in (13) and (14) and $Q_{D,R}$ be RDAQ described in Alg. 10 and 11. Then, for $\mu d \in [d]$, $\mathcal{X} = \mathcal{Y} = \mathcal{B}$, and every $\Delta > 0$, we have*

$$\alpha(Q_{\text{wz},u}; \Delta) \leq \frac{\alpha(Q_{D,R}; \Delta)}{\mu} \quad \text{and} \quad \beta(Q_{\text{wz},u}; \Delta) = 0.$$

Furthermore, the output of quantizer $Q_{\text{wz},u}$ can be described in $\mu d(h + \log h)$ bits.

We are now equipped to prove our second main result. Our protocol π_u^* uses $Q_{\text{wz},u}$ for each client as described in Section 2 and forms the estimate \hat{x} as in (4). Unlike for the known Δ setting, we now use the same parameters for $Q_{\text{wz},u}$ for all clients, given by

$$\mu d = \left\lfloor \frac{r}{h + \log h} \right\rfloor. \quad (15)$$

Theorem 4.4. *For $d \geq r \geq 2(h + \log h)$ and h given in (12), the r -bit protocol π_u^* with parameters as set in (15) satisfies*

$$\text{MSE}(\pi_u^*, \mathbf{x}, \mathbf{y}) \leq (128\sqrt{3}(1 + \ln^*(d/6))) \left(\sum_{i \in [n]} \frac{\Delta_i}{n} \cdot \frac{d}{nr} \right),$$

for all \mathbf{x}, \mathbf{y} satisfying (2), for every $\Delta = (\Delta_1, \dots, \Delta_n)$.

Proof. Denote by \hat{x} the output of the protocol. Then, by Lemmas 2.1 and Lemma 4.3, we get

$$\begin{aligned}\mathbb{E} [\|\hat{x} - \bar{x}\|_2^2] &\leq \frac{1}{n^2\mu} \sum_{i=1}^n \alpha(Q_{D,R}; \Delta_i) \\ &\leq \frac{16\sqrt{3}}{n^2\mu} \sum_{i=1}^n \Delta_i,\end{aligned}$$

where the previous inequality is by Lemma 4.2. The proof is completed by using $\mu \geq \frac{r}{2d(h+\log h)} \geq \frac{r}{4dh}$, which follows from (15) and the assumption that $r \geq 2(h + \log h)$. \square

The Wyner-Ziv estimator π_u^* is universal in Δ : it operates without the knowledge of the distance between the input and the side information and yet gets MSE depending on Δ . Moreover, it can be efficiently implemented as both the encoding and the decoding procedures have nearly linear time complexity of $O(d \log d)$.

5 Application: Communication constrained distributed optimization

We consider the problem of minimizing an unknown convex function $f: \mathcal{X} \rightarrow \mathbb{R}$ over its domain $\mathcal{X} \subset \mathbb{R}^d$ using the set of n clients who have access to independent noisy gradients of the function. In particular, the server runs an optimization algorithm, which is not directly given access to the function but can get n different gradient estimates of the function at various points of its choice. This class of optimization algorithms includes various descent algorithms, which provide close to optimal convergence rate within the class and are appealing in practice due to their distributed nature.

Owing to our setup, the gradient estimates supplied by the n clients must pass through r -bit quantizers, chosen from a fixed set of quantizers \mathcal{Q}_r ¹¹, and the optimization algorithm **A** only has access to the quantized outputs.

Our objective is to select quantizers $Q_{i,t}$, $\forall i \in [n], t \in [T]$, and an optimization algorithm **A** to guarantee the minimum worst-case optimization error defined below. In our setting, we allow for *adaptive gradient processing*, whereby, the quantizer $Q_{i,t}$ selected in t th iteration may depend on all the previous quantized outputs. Specifically, denoting by $C_{i,t}$ the i th client's quantized output at time t , which takes values in the output alphabet \mathbb{R}^d , the *adaptive quantizer selection strategy* $S := (S_1, \dots, S_T)$ over T iterations consists of mappings $S_t: \mathbb{R}^{d \times n \times (t-1)} \rightarrow \mathcal{Q}_r^n$ that take $\{C_{i,t'}\}_{i \in [n], t' \in [t-1]}$ as input and outputs a tuple of n quantizers $\{Q_{i,t}\}_{i \in [n]} \in \mathcal{Q}_r^n$. We write $\mathcal{S}_{\mathcal{Q}_r, T}$ for the collection of all such quantizer selection strategies. The entire framework can be summarized as follows:

1. At iteration t , the first-order optimization algorithm **A** makes a query for point x_t to clients $\mathbf{Cl}_1, \dots, \mathbf{Cl}_n$.
2. Upon receiving the point $x_t \in \mathcal{X}$, the client \mathbf{Cl}_i , $i \in [n]$, outputs $\hat{g}_i(x_t)$, an unbiased estimate of $\nabla f(x_t)$.

¹¹The set of r -bit quantizers \mathcal{Q}_r is used to model the communication constraints in a distributed setting.

3. The gradient estimate $\hat{g}_i(x_t)$ is passed through a quantizer $Q_{i,t} \in \mathcal{Q}_r$ chosen based on strategy S , and the output $Y_{i,t}$ is observed by the first-order optimization algorithm \mathbf{A} . The algorithm then uses all the messages $\{C_{i,t'}(x_{t'})\}_{i \in [n], t' \in [t]}$ to further update x_t to x_{t+1} .

Denote by \mathbf{C} the collection of n clients $(\mathbf{C}_1, \dots, \mathbf{C}_n)$. Let \mathcal{A}_T be the set of all first-order optimization algorithms that make T queries to \mathbf{C} and for the t th query x_t , get back the outputs $\{Y_{i,t}\}_{i \in [n]}$. We measure the performance of an optimization protocol \mathbf{A} and a quantizer selection strategy S for a given function f and clients $\mathbf{C}_i, i \in [n]$, using the metric $\mathcal{E}(f, \mathbf{C}, \mathbf{A}, S)$ defined as

$$\mathcal{E}(f, \mathbf{C}, \mathbf{A}, S) = \mathbb{E} \left[f(\bar{x}_T) - \min_{x \in \mathcal{X}} f(x) \right],$$

where $\bar{x}_T := \frac{1}{T} \sum_{t \in [T]} x_t$ and the expectation is over the randomness in \bar{x}_T .

For a set of various function and client pairs above, denoted by \mathcal{O} , the set of r -bit quantizers \mathcal{Q}_r and the number of iterations T , we define the *minimax optimization error* as

$$\mathcal{E}^*(\mathcal{X}, \mathcal{O}, T, \mathcal{Q}_r) = \inf_{\mathbf{A} \in \mathcal{A}_T} \inf_{S \in \mathcal{S}_{\mathcal{Q}_r, T}} \sup_{(f, \mathbf{C}) \in \mathcal{O}} \mathcal{E}(f, \mathbf{C}, \mathbf{A}, S).$$

We now define the class of functions and state the assumptions related to the clients accessible to the algorithm \mathbf{A} .

Convex and smooth function family Throughout, we restrict ourselves to convex and L -smooth functions over $\mathcal{X} \subset \mathbb{R}^d$, i.e., functions satisfying, $\forall \lambda \in [0, 1], \forall x, y \in \mathbb{R}^d$,

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y), \quad (16)$$

$$\|\nabla f(x) - \nabla f(y)\|_2 \leq L\|x - y\|_2, \quad (17)$$

where $\nabla f(x) \in \mathbb{R}^d$ denotes the gradient of f at input x .

Stochastic gradients We assume that the output $\hat{g}_i(x)$ by client $\mathbf{C}_i, 1 \leq i \leq n$, when a point $x \in \mathcal{X}$ is queried satisfies the following conditions:

$$\mathbb{E}[\hat{g}_i(x) \mid x] = \nabla f(x), \quad (\text{unbiased estimates}) \quad (18)$$

$$\|\hat{g}_i(x) - \nabla f(x)\|_2^2 \leq \sigma^2, \quad (\text{maximum deviation bound}) \quad (19)$$

$$\|\hat{g}_i(x)\|_2^2 \leq B^2. \quad (\text{a.s. bounded estimate}) \quad (20)$$

Assumption (18) is standard in stochastic optimization literature (*cf.* [45], [44], [13]). However, it is enough to assume a bound on the variance of stochastic gradients instead of (19) to prove convergence guarantees for smooth stochastic optimization without any communication constraints. The stronger assumption made here is to aid a much tighter analysis under communication constraints. In Section 5.5, we provide a scheme which can operate under the standard variance bound.

Denote by \mathcal{O}_{sc} the set of tuples of function and n clients, (f, \mathbf{C}) , satisfying (16), (17), (18), (19) and (20).

5.1 Lower bound

The following bound will serve as a basic benchmark for our problem. Let $D > 0$ and $\mathbb{X}_2(D) := \{\mathcal{X} \subseteq \mathbb{R}^d : \max_{x,y \in \mathcal{X}} \|x - y\|_2 \leq D\}$ be the collection of subsets of \mathbb{R}^d whose ℓ_2 diameter is at most D .

Theorem 5.1. *There exists an absolute constant $0 \leq c_0 \leq 1$ such that for $r \leq d$ and $T \geq d/(6nr)$,*

$$\sup_{\mathcal{X} \in \mathbb{X}_2(D)} \mathcal{E}^*(\mathcal{X}, \mathcal{O}_{\text{sc}}, T, \mathcal{Q}_r) \geq \frac{c_0 D \sigma}{\sqrt{nT}} \cdot \sqrt{\frac{d}{r}}.$$

5.2 A general convergence bound

We present a general convergence bound based on a non-adaptive channel strategy. In particular, we fix same quantization process in every iteration, and the quantized outputs $\{C_{i,t}\}_{i \in [n]}$ are passed through a mapping¹² $\mathcal{M} : \mathbb{R}^{n \times d} \rightarrow \mathbb{R}^d$ in order to update the query.

We use PSGD as the first-order optimization algorithm; the overall optimization procedure is described in Algorithm 12. PSGD proceeds as SGD, with the additional projection step where it projects the updates back to domain \mathcal{X} using the map $\Gamma_{\mathcal{X}}(y) := \min_{x \in \mathcal{X}} \|x - y\|$, $\forall y \in \mathbb{R}^d$.

1: **for** $t = 0$ to $T - 1$ **do**
 2: $x_{t+1} = \Gamma_{\mathcal{X}}(x_t - \eta_t \mathcal{M}(C_{1,t}, \dots, C_{n,t}))$
 3: **Output** $\bar{x}_T = \frac{1}{T} \sum_{t=1}^T x_t$

Algorithm 12: PSGD using clients \mathcal{C}

The convergence rate of Algorithm 12 is controlled by the worst-case L_2 -norm $\alpha'(\mathcal{M})$ and the worst-case bias $\beta'(\mathcal{M})$ defined as

$$\alpha'(\mathcal{M}) := \sup_{\substack{\{\forall x, i \in [n], \hat{g}_i \in \mathbb{R}^d: \\ \|\hat{g}_i - \nabla f(x)\|^2 \leq \sigma^2\}}} \sqrt{\mathbb{E}[\|\mathcal{M}(C^n) - \nabla f(x)\|^2]}, \quad (21)$$

$$\beta'(\mathcal{M}) := \sup_{\substack{\{\forall x, i \in [n], \hat{g}_i \in \mathbb{R}^d: \\ \|\hat{g}_i - \nabla f(x)\|^2 \leq \sigma^2\}}} \|\mathbb{E}[(\mathcal{M}(C^n) - \nabla f(x))]\|, \quad (22)$$

where $C^n = (C_1, \dots, C_n)$ is the communication received at the server. Using a slight modification of the standard proof of convergence for PSGD in [13, Theorem 6.3], we can derive the following lemma.

Lemma 5.2. *For any mapping \mathcal{M} and set of quantizers $\{Q_i\}_{i \in [n]}$ defined above, the output \bar{x}_T of optimization algorithm given in Algorithm 12 satisfies*

$$\sup_{(f, \mathcal{C}) \in \mathcal{O}} \mathcal{E}(f, \mathcal{C}, \mathbf{A}, S) \leq \frac{\sqrt{2}\alpha'(\mathcal{M})D}{\sqrt{T}} + \beta'(\mathcal{M}) \left(D + \frac{DB}{\alpha'(\mathcal{M})\sqrt{2T}} \right) + \frac{LD^2}{2T},$$

with the learning rate $\eta_t = \frac{1}{L + \frac{\alpha'(\mathcal{M})\sqrt{2T}}{D}}, \forall t \in [T]$.

¹²For instance, averaging the quantized outputs at the server can possibly be one such mapping.

5.3 Baseline scheme: Parallel SGD

We begin by presenting the convergence result for the baseline scheme in our setup: the `Parallel SGD` algorithm. In `Parallel SGD`, all clients compress their stochastic gradient estimates to r bits using an efficient quantizer for the Euclidean ball and send it to the server, which then takes the average of the quantized gradients for the projected gradient descent step. We choose subsampled RATQ ([43]) for this efficient quantizer. We denote by Q_{RATQ} the subsampled version of RATQ using r bits, which is described in [43, Section 3.5]. After receiving the quantized outputs $C_{i,t} = Q_{\text{RATQ}}(\hat{g}_i(x_t)) \forall i \in [n]$, from all the n clients, the server takes the mapping \mathcal{M} to be the average of these outputs, i.e.,

$$\mathcal{M}(\bar{C}_t) = \frac{1}{n} \sum_{i=1}^n Q_{\text{RATQ}}(\hat{g}_i(x_t)). \quad (23)$$

Theorem 5.3. *Let S be the quantizer selection strategy which fixes the quantizer to be Q_{RATQ} for all clients at all iterations. Let \mathbf{A} be the optimization algorithm described in Algorithm 12 where \mathcal{M} as described in (23) is used to make the PSGD step after the t -th query and the learning rate $\eta_t = \frac{1}{L + \frac{\alpha'(\mathcal{M})\sqrt{2T}}{D}}$, where $\alpha'(\mathcal{M}) = c_0 \sqrt{\frac{\sigma^2}{n} + \frac{c_2 dB^2 \log \log^* d}{nr}}$ for some positive universal constant c_0 . Then, for positive universal constants c_1 and c_2 and r such that $d \geq r \geq c_1 \log \log^* d$, we have*

$$\mathcal{E}(f, \mathbf{C}, \mathbf{A}, S) \leq \frac{c_2 D}{\sqrt{nT}} \sqrt{\sigma^2 + \frac{c_2 dB^2 \log \log^* d}{r}} + \frac{LD^2}{2T}.$$

We note that the term $\frac{dB^2 \log \log^* d}{r}$ illustrates the slowdown in convergence due to quantization error. This is nearly the best rate which can be achieved when one uses r -bit quantizers without any side information¹³. Note that in cases in which B is large relative to σ^2 , the slowdown due to this term can be significant, and the algorithm maybe far away from our lower bound in Theorem 5.1.

5.4 WZ-SGD: An almost optimal algorithm for distributed optimization

We now present our main algorithm : `WZ-SGD`. `WZ-SGD` uses our first Wyner-Ziv estimator (see Section 3.3) based on subsampled RMQ as a subroutine to form much more accurate gradient estimates compared to those formed in `ParallelSGD`. As a result of this, `WZ-SGD` significantly improves over the convergence rate of Theorem 5.3 and relegates the dependence of convergence rate on B to only second order terms.

At each iteration t , `WZ-SGD` uses the clients in \mathbf{C}_1 to form the side information estimate Z_t at the server and then uses the clients in \mathbf{C}_2 to estimate the gradient for performing the descent step, where¹⁴ $\mathbf{C}_1 := \{\mathbf{C}_{1,1}, \dots, \mathbf{C}_{1,n/2}\}, \mathbf{C}_2 := \mathbf{C} \setminus \mathbf{C}_1$.

The side information estimate Y_t . The side information is formed as follows. Under the r -bit communication constraint, we divide the coordinates into blocks of dimension r_1 , where

¹³Similar convergence bounds (upto $\log \log d$ factor) for parallel SGD can be achieved by using subsampled version of rotated quantizer in [55] or the subsampled version of uniform quantizer after preprocessing due to Kashin's representation (cf. [28], [39]).

¹⁴For simplicity, we assume that $n/2$ and d/r_1 are integers such that d/r_1 divides $n/2$.

```

1: for Clients  $i \in [n]$  do    ▷ Setting quantizers
2:   if  $i \in \mathbf{C}_1$  then  $Q_i = Q_u$ 
3:   else  $Q_i = Q_{\text{wz},i}$ 
4: Initialize  $x_0 \in \mathcal{X}$ 
5: for  $t = 0$  to  $T - 1$  do
6:   for Server do
7:     Broadcast  $x_t$  to clients
8:   for Clients  $i \in [n]$  do    ▷ Encoding
9:     Compute  $\hat{g}_i(x_t)$ 
10:    Send  $Q_i^e(\hat{g}_i(x_t))$  to server
11:  for Server do    ▷ Decoding
12:    for  $i \in \mathbf{C}_1$  do
13:       $Q_i(\hat{g}_i(x_t)) = Q_i^d(Q_i^e(\hat{g}_i(x_t)))$ 
14:       $Y_t = \frac{2}{n} \sum_{i \in \mathbf{C}_1} Q_i(\hat{g}_i(x_t))$  ▷ Side information
15:      for  $i \in \mathbf{C}_2$  do
16:         $Q_i(\hat{g}_i(x_t), Y_t) = Q_i^d(Q_i^e(\hat{g}_i(x_t)), Y_t)$ 
17:       $x_{t+1} = \Gamma_{\mathcal{X}}(x_t - \eta_t \cdot \frac{2}{n} \sum_{i \in \mathbf{C}_2} Q_i(\hat{g}_i(x_t), Y_t))$ 
18: At Server Output:  $\bar{x}_T = \frac{1}{T} \sum_{t \in [T]} x_t$ 

```

Algorithm 13: WZ-SGD algorithm

$r_1 := r / \log \ell_1$, and $\log \ell_1$ denotes the precision bits used by clients to represent each coordinate in the assigned block. This way we have d/r_1 blocks. We also equally partition the set \mathbf{C}_1 into d/r_1 groups. Further, we assign every block of r_1 coordinates to every other distinct group of $\frac{n/2}{d/r_1}$ clients. To quantize the coordinates within any block, the group of clients assigned to that block will use a coordinate-wise uniform quantizer (CUQ). CUQ is an unbiased, uniform quantizer that has appeared recently in many works on gradient quantization. We denote by $Q_u : [-B, B] \rightarrow \{-B + 2B \cdot (i-1)/(\ell_1-1) : i \in [\ell_1]\}$ the ℓ_1 -level CUQ quantizer. For a scalar input $x \in [-B, B]$,

$$Q_u(x) = \begin{cases} \left\lceil \frac{x(\ell_1-1)}{2B} \right\rceil \cdot \frac{2B}{\ell_1-1}, & \text{w.p. } \frac{x - \lfloor \frac{x(\ell_1-1)}{2B} \rfloor}{\frac{2B}{\ell_1-1}}, \\ \lfloor \frac{x(\ell_1-1)}{2B} \rfloor \cdot \frac{2B}{\ell_1-1}, & \text{w.p. } \frac{\lceil \frac{x(\ell_1-1)}{2B} \rceil - x}{\frac{2B}{\ell_1-1}}. \end{cases} \quad (24)$$

Each client uses an ℓ_1 -level CUQ to quantize the associated block of coordinates separately. Thus, the overall communication by each client is $r_1 \cdot \log \ell_1 = r$ and satisfies the communication constraint.

For each block, we then form the side information by taking the average of the quantized outputs from all its associated clients. Denote by Y_t the side information formed at the server by using the clients in \mathbf{C}_1 at iteration t . Then, from the description of our scheme, for all coordinates $i \in \{r_1(j-1) + 1, \dots, r_1 j\}$ and for all $j \in [d/r_1]$ we have

$$Y_t(i) = \frac{2d}{nr_1} \sum_{k \in \mathcal{S}_j} Q_u(\hat{g}_k(x_t)(i)),$$

where \mathcal{S}_j denotes the set of $\frac{nr_1}{2d}$ clients assigned to form the side information for the coordinates

$\{r_1(j-1)+1, \dots, r_1 j\}$, i.e.,

$$\mathcal{S}_j = \{\mathbf{C1}_{(nr_1/(2d)) \cdot (j-1) + 1}, \dots, \mathbf{C1}_{(nr_1/(2d)) \cdot j}\}. \quad (25)$$

We remark that to decode each quantized gradient estimate sent by clients in \mathcal{C}_2 , we will use Y_t as side information. However, Y_t will not be used as is but a version which is rotated¹⁵ using a random matrix (8) will be used.

The Wyner-Ziv gradient estimate Q_{WZ} . We use the clients in \mathcal{C}_2 to form the actual gradient estimate. The clients encode the stochastic gradients using a subsampled RMQ quantizer (see Section 3.3 for details). Therefore, for stochastic gradient $\hat{g}_j(x_t)$, the output encoded by client $\mathbf{C1}_j$ using subsampled RMQ is described as follows:

$$Q_{\text{WZ},j}^e(\hat{g}_j(x_t)) = \{Q_{\mathbf{M}}^e(R_j \hat{g}_j(x_t)(i)) : i \in \mathcal{D}_j\}.$$

At the server, the communication for all $\mathbf{C1}_j \in \mathcal{C}_2$ is decoded as follows:

$$Q_{\text{WZ},j}(\hat{g}_j(x_t), Y_t) = R_j^{-1} \left(\frac{d}{r_2} \sum_{i \in \mathcal{D}_j} (\tilde{g}_j - R_j Y_t(i)) e_i + R_j Y_t \right)$$

where $\tilde{g}_j(i) = Q_{\mathbf{M}}(R_j \hat{g}_j(x_t)(i), R_j Y_t(i))$. Finally, the server averages over all the quantized gradient estimates of clients in \mathcal{C}_2 to get (see, line 17 in Algorithm 13)

$$\mathcal{M}(Q_{\text{WZ},1}, \dots, Q_{\text{WZ},n}) = \frac{2}{n} \sum_{j=n/2+1}^n Q_{\text{WZ},j}(\hat{g}_j(x_t), Y_t) \quad (26)$$

Next, we present the convergence rate of the proposed WZ-SGD algorithm for communication constrained distributed optimization.

Theorem 5.4. *Let S be the communication protocol which uses the CUQ quantizer for clients \mathcal{C}_1 and the subsampled RMQ quantizer for clients in \mathcal{C}_2 . Let \mathbf{A} be the optimization algorithm described in Algorithm 13 with the learning rate $\eta_t = \frac{1}{L + \frac{\alpha'(\mathcal{M})\sqrt{2T}}{D}}$, where $\alpha'(\mathcal{M}) = c_0 \sqrt{\frac{d\sigma^2 \log \log nT}{nr}}$ for some positive universal constant c_0 . Then, for positive universal constants c_1, c_2 , and c_3 and r, n such that $d \geq r \geq c_1 \max\{\log \log nT, \log(B/\sigma)\}$ and $nr \geq c_2 d^2 \log(B/\sigma)$, we have*

$$\mathcal{E}(f, \mathcal{C}, \mathbf{A}, S) \leq \frac{c_3 D \sigma}{\sqrt{nT}} \cdot \sqrt{\frac{d \log \log nT}{r}} + \frac{LD^2}{2T}.$$

Remark 5. The condition on nr is needed to remove any B dependence from the MSE upper bound.

Thus, in the setting where the number of clients n is large, we match the lower bound in Theorem 5.1 upto a $\log \log nT$ factor.

5.5 UWZ-SGD: A universal Wyner-Ziv algorithm for distributed optimization

We now relax the almost sure (19) assumption on the gradients estimated by clients and present an *universal* algorithm UWZ-SGD, where the compression at the clients doesn't need the knowledge

¹⁵For decoding each quantized gradient sent by clients in \mathcal{C}_2 , Y_t will be rotated using independent and identical versions of matrix R .

```

1: for Clients  $i \in [n]$  do    ▷ Setting quantizers
2:   if  $i \in \mathcal{C}_1$  then  $Q_i = Q_{\text{RATQ}}$ 
3:   else  $Q_i = Q_{\text{WZ,u}}$ 
4: Initialize  $x_0 \in \mathcal{X}$ 
5: for  $t = 0$  to  $T - 1$  do
6:   for Server do
7:     Broadcast  $x_t$  to clients
8:   for Clients  $i \in [n]$  do    ▷ Encoding
9:     Compute  $\hat{g}_i(x_t)$ 
10:    Send  $Q_i^e(\hat{g}_i(x_t))$  to server
11:  for Server do    ▷ Decoding
12:    for  $i \in \mathcal{C}_1$  do
13:       $Q_i(\hat{g}_i(x_t)) = Q_i^d(Q_i^e(\hat{g}_i(x_t)))$ 
14:       $Y_t = \frac{2}{n} \sum_{i \in \mathcal{C}_1} Q_i(\hat{g}_i(x_t))$  ▷ Side information
15:      for  $i \in \mathcal{C}_2$  do
16:         $Q_i(\hat{g}_i(x_t), Y_t) = Q_i^d(Q_i^e(\hat{g}_i(x_t)), Y_t)$ 
17:       $x_{t+1} = \Gamma_{\mathcal{X}}(x_t - \eta_t \cdot \frac{2}{n} \sum_{i \in \mathcal{C}_2} Q_i(\hat{g}_i(x_t), Y_t))$ 
18: At Server Output:  $\bar{x}_T = \frac{1}{T} \sum_{t \in [T]} x_t$ 

```

Algorithm 14: UWZ-SGD algorithm

of σ and only the server needs to know σ to set the learning rate in Algorithm 12. Specifically, we assume that for all clients $i \in [n]$,

$$\mathbb{E} [\|\hat{g}_i(x) - \nabla f(x)\|^2] \leq \sigma^2. \quad (\text{m.s. deviation bound}) \quad (27)$$

The other assumptions (18) and (20) about the estimated gradients still hold¹⁶. We show how the dependence of B in the naive scheme, presented in Theorem 5.3, can be reduced using subsampled RDAQ.

At every iteration, the client indexed by \mathcal{C}_1 use subsampled RATQ to compress their gradient estimates. The side information is then formed by taking sample average of the decoded estimates, similar to (25) (see line 14, in Algorithm 14).

On the other hand, the clients in \mathcal{C}_2 use the subsampled RDAQ quantizer $Q_{\text{WZ,u}}$ from section 4.4. Note that the subsampled RDAQ decoder (14) uses the side information constructed by \mathcal{C}_1 . Finally, the server takes the sample average of the decoded values estimated by the \mathcal{C}_2 (see, line 17 in Algorithm 14) to form the mapping \mathcal{M} .

Theorem 5.5. *Let S be the communication protocol which uses the subsampled RATQ quantizer for clients \mathcal{C}_1 and the subsampled RDAQ quantizer for clients in \mathcal{C}_2 . Let \mathbf{A} be the optimization algorithm described in Algorithm 14 with the learning rate $\eta_t = \frac{1}{L + \frac{\alpha'(\mathcal{M})\sqrt{2T}}{B}}$, where $\alpha'(\mathcal{M}) =$*

$\sqrt{\frac{2\sigma^2}{n} + \frac{2\rho(B,\sigma,r,n)}{n}}$ with $\rho = \sqrt{\sigma^2 + \frac{2\sigma^2}{n} + \frac{2dB^2}{n \left(\frac{r}{3 + \lceil \log(1 + \ln^(d/3)) \rceil} - 1 \right)}} \frac{16\sqrt{3}dB}{\lceil h + \log h \rceil - 1}$ and $h = 1 + \ln^*(d/6)$. Further, suppose that the gradient estimated by all the clients satisfy the assumptions (18), (27) and*

¹⁶Note that the lower bound in Theorem 5.1 under the almost sure assumption (19) holds for the relaxed mean-squared assumption (27) too.

(20). Then, for $d \geq r \geq \max\{h + \log h, 3 + \lceil \log(1 + \ln^*(d/3)) \rceil\}$, we have

$$\mathcal{E}(f, \mathbf{c}, \mathbf{A}, S) \leq \frac{2D}{\sqrt{nT}} \sqrt{\sigma^2 + \rho(B, \sigma, r, n)} + \frac{LD^2}{2T}.$$

Remark 6. We remark that under the relaxed assumption of mean-square bounded deviation in (27), for $nr \geq (B^2/\sigma^2)d \log(1 + \ln^*(d/3))$, the slowdown in the convergence rate is illustrated by $\rho \approx \frac{16\sqrt{3}dB\sigma \ln^* d}{r}$, and the universal scheme surpasses the performance of parallel SGD presented in Section 5.3.

We end this section by pointing out limitations of a natural scheme for distributed optimization.

Remark 7 (Limitations of Centering Based Scheme). We note that our framework allows for quantization schemes where previously quantized gradients are used for gradient compression at the current iteration. For instance, we can use average of the compressed gradients at the previous iteration to center the current compression. That is, the server broadcasts the average to all the clients and the clients only need to compress the difference between the current stochastic gradient and this communicated average.

If the query points x_{t-1} and x_t do not deviate by much, then such type of compression schemes which are centered around the average of previous quantized gradients may turn out to be very efficient. Also, note that the typical learning rate for smooth optimization is $O(\sqrt{1/T})$, which means that the difference between the points x_t and x_{t-1} is not very large. Moreover, the smoothness assumption (17) allows to control the deviation between the true gradients at successive iterations in terms of the points queried at the two iterations. All this hints at the fact that such a scheme where each client uses optimal quantizers for quantizing the difference vector without any side-information may turn out to be optimal. But note that for a very large value of smoothness constant, $L \geq \sigma\sqrt{T}/D$, even with small deviation between successive query points, the deviation between the gradients will be large. This would in turn lead to variance of the quantized gradients having a dependence on the maximum gradient norm $\sqrt{B^2 - \sigma^2}$, which would in turn lead to the leading term, in terms of n and T , in convergence rate depending on $\sqrt{B^2 - \sigma^2}$.

6 The Gaussian Wyner-Ziv problem

Consider the random vectors X and Y , where the coordinates $\{X(i), Y(i)\}_{i=1}^d$ form an i.i.d. sequence. Furthermore, for all $i \in [d]$, let

$$X(i) = Y(i) + Z(i),$$

where $Y(i)$ and $Z(i)$ are independent and zero-mean Gaussian random variables with variances σ_y^2 and σ_z^2 , respectively. The encoder has access to the sequence $X = \{X(i)\}_{i=1}^d$, which it quantizes and sends to the decoder. The decoder, on the other hand, has access to Y (note that encoder does not have access to Y) and can use it to decode X . A pair (R, D) of non-negative numbers is an achievable rate-distortion pair if we can find a quantizer Q_d of precision dR and with mean square error $\mathbb{E}[\|Q_d(X, Y) - X\|_2^2] \leq dD$. For $D \geq 0$, denote by $R(D)$ the infimum over all R such that (R, D) constitute an achievable rate-distortion pair for all d sufficiently large. From¹⁷ [61], $R(D)$

¹⁷The model considered in [61] and perhaps the more popular Wyner-Ziv model is $Y = X + Z$. Nevertheless, through MMSE rescaling this model can be converted to $X = Y' + Z'$ (see, for instance, [36]).

can be characterized as follow:

$$R(D) = \begin{cases} \frac{1}{2} \log \frac{\sigma_z^2}{D} & \text{if } D \leq \sigma_z^2, \\ 0 & \text{if } D > \sigma_z^2. \end{cases}$$

Several constructions that involve computational heavy methods such as error correcting codes and lattice encoding attain the rate-distortion function, asymptotically for large d . In this section, we show that modulo quantizer with parameters set appropriately attains a rate very close to the rate-distortion function $R(D)$. Moreover, we will show that this rate can be achieved for arbitrary Y and Z , as long as Z is a zero mean subgaussian random variable with variance factor σ_z^2 . Our proposed quantizer $Q_d(X, Y)$ uses the modulo quantizer to quantize $X(i)$ with side information $Y(i)$ at the decoder and the parameter k, Δ' set as follows:

$$\begin{aligned} \delta &= \sqrt{D/308}, \quad \log k = \left\lceil \log \left(2 + \sqrt{\frac{24\sigma_z^2}{D} \ln \frac{308\sigma_z^2}{D}} \right) \right\rceil \\ \Delta' &= \sqrt{6(\sigma_z^2) \ln(\sigma_z/\delta)}, \quad \varepsilon = 2\Delta'/(k-2), \end{aligned} \quad (28)$$

Theorem 6.1. *Consider random vectors X, Y in \mathbb{R}^d with $X(i) = Y(i) + Z(i)$ and $Z(i)$ independent of $Y(i)$ being a centered subgaussian random variable with variance factor of σ_z^2 , for all coordinates $i \in \{1, \dots, d\}$. Then, for $D \leq (\sigma_z^2/308)$, the quantizer $Q_d(X, Y)$ described above has MSE less than dD and has rate R satisfying*

$$R \leq \frac{1}{2} \log \frac{\sigma_z^2}{D} + O\left(\log \log \frac{\sigma_z^2}{D}\right).$$

7 The high-precision regime

7.1 RMQ in the high-precision regime.

For the known Δ setting, our quantizer RMQ described in Alg. 4 and 5 remains valid even for $r > d$. We will assume $r = md$ for integer $m \geq 2$. For each client i , we set

$$\delta = \frac{\Delta_i}{n^{\frac{1}{2}(2^{r/d} - 2)}}, \quad \log k = \frac{r}{d}, \quad \Delta' = \sqrt{6(\Delta_i^2/d) \ln \Delta_i/\delta}, \quad \varepsilon = \frac{2\Delta'}{k-2}. \quad (29)$$

The performance of protocol π_k^* using RMQ with parameters set as in (29) for each client can be characterized as follows.

Theorem 7.1. *For a fixed $\Delta = (\Delta_1, \dots, \Delta_n)$ and $r = md$ for integer $m \geq 2$, the protocol π_k^* with parameters set as in (29) satisfies*

$$\text{MSE}(\pi_k^*, \mathbf{x}, \mathbf{y}) \leq \left(12 \ln n + \frac{24r}{d} + 154/n + 166 \right) \left(\sum_{i \in [n]} \frac{\Delta_i^2}{n} \cdot \frac{1}{n(2^{r/d} - 2)^2} \right),$$

for all \mathbf{x}, \mathbf{y} satisfying (2).

Proof. Denoting by Q_i the quantizer $Q_{M,R}$ with parameters set for client i , by Lemmas 2.1 and 3.2, we get

$$\mathbb{E} [\|\hat{x} - \bar{x}\|_2^2] \leq \sum_{i=1}^n \frac{\alpha(Q_i; \Delta_i)}{n^2} + \sum_{i=1}^n \frac{\beta(Q_i; \Delta_i)}{n}$$

Further, since $k \geq 4$ holds when $r \geq 2d$ for our choice of parameters, by using Lemma 3.2 and substituting $\delta^2 = \Delta_i^2/n(2^{r/d} - 2)^2$, we get

$$\begin{aligned} \alpha(Q_i; \Delta_i) &\leq \frac{12\Delta_i^2 \ln(n(2^{r/d} - 2)^2)}{(2^{r/d} - 2)^2} + \frac{154\Delta_i^2}{n(2^{r/d} - 2)^2}, \\ \beta(Q_i; \Delta_i) &\leq \frac{154\Delta_i^2}{n(2^{r/d} - 2)^2}. \end{aligned}$$

which with the previous bound gives

$$\mathbb{E} [\|\hat{x} - \bar{x}\|_2^2] \leq \left(12 \ln n + \frac{24r}{d} + \frac{154}{n} + 154\right) \sum_{i=1}^n \frac{\Delta_i^2}{n^2(2^{r/d} - 2)^2},$$

where use the inequality $\ln x \leq x, \forall x \geq 0$, to bound $\ln(2^{r/d} - 2)^2/(2^{r/d} - 2)^2$ by 1. □

Remark 8. Similar to Remark 4, we note that using MQ for each coordinate without rotating (or even with rotation using R as above) with $\Delta' = \Delta_i$ yields MSE less than

$$O\left(\sum_{i=1}^n \frac{\Delta_i^2}{n} \cdot \frac{d}{n2^{2r/d}}\right),$$

for $r \geq d$. Thus our approach above allows us to remove the d factor at the cost of a (milder for large d) $\log n + r/d$ factor.

7.2 Boosted RDAQ: RDAQ in the high-precision regime.

Moving to the unknown Δ setting, we describe an update to RDAQ described in Alg. 10 and 11 for the high-precision setting. For brevity, we denote by $m := r/d$ the number of bits per dimension. A straight-forward scheme to make use of the high precision is to independently implement the RDAQ quantizer approximately $\lfloor m/\ln^* d \rfloor$ times and use the average of the quantized estimates as the final estimate. We will see that the MSE incurred by such an estimator is $O(\Delta \ln^* d/m)$. We will show that this naive implementation can be significantly improved and an exponential decay in MSE with respect to m can be achieved.

We boost RDAQs performance as follows. Simply speaking, instead of sending the bits produced by multiple instances of the encoder of RDAQ, we send the “type” of each sequence. A similar idea appeared in [42] for the case without any side information. At the encoding stage of RDAQ given in Alg. 10 and 11, after random rotation and computing z in Steps 1 to 3 of Alg. 10, we repeat Step 4 N times with independent randomness each time and store only the total number of ones seen for each coordinate i and scale j . Specifically, let $U_t(i, j)$ be an independent uniform

random variable in $[-M_j, M_j]$, for all $i \in [d], j \in [h]_0$, and $t \in [N]$, which are generated using public randomness between the encoder and the decoder. Using this randomness, we compute $\tilde{x}_{j,t} = \sum_{i=1}^d \mathbb{1}_{\{U_t(i,j) \leq x_{R(i)}\}} e_i$ for all $j \in [h]_0$. Then, instead of storing $\tilde{x}_{j,t}$ for each j and t , we store the sum $\sum_{t=1}^N \tilde{x}_{j,t}$ for each $j \in [h]_0$. Since each coordinate of the sum can be stored in $\log(N+1)$ bits, the new encoder's output can be stored in $d(h \log(N+1) + \log h)$. Thus, we can implement this scheme by using $m = (h \log(N+1) + \log h)$ bits per dimension.

At the decoding stage, we rotate y and compute z^* in precisely the same manner as done in Steps 1 to 3 of the decoding Alg. 11 of RDAQ. Then, using the encoded input received, the side-information y , the same random variables $U_t(i, j)$ and random matrix R used by the encoder, the final estimate $Q(x)$ is

$$Q(x) = R^{-1} \left(\frac{1}{N} \cdot \sum_{i \in [d]} \sum_{t \in [N]} (B_{i,Rx}^t - B_{i,Ry}^t) e_i + Ry \right), \quad (30)$$

where $B_{i,v}^t = \mathbb{1}_{\{U_t(i,z^*(i)) \leq v(i)\}}$ for v in \mathbb{R}^d .

The result below characterizes the performance of our quantizer Boosted RDAQ Q .

Lemma 7.2. *Let Q be Boosted RDAQ described above. Then, we have for $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$ and every $\Delta > 0$, we have*

$$\alpha_u(Q; \Delta) \leq \frac{16\sqrt{3}\Delta}{N} \quad \text{and} \quad \beta_u(Q; \Delta) = 0.$$

Furthermore, the output of the quantizer can be described in $d(h \log N + \log h)$ bits.

Thus, when we have a total precision budget of $r = dm$ bits using the Boosted RDAQ algorithm with number of repetitions $N = 2^{\lfloor (m - \log h)/h \rfloor} - 1$, we get an exponential decay in MSE with respect to m .

We consider the protocol π_u^* that uses the Q above for each client with M_j and h set as in (12), i.e., with

$$N = 2^{\lfloor (m - \log h)/h \rfloor} - 1, \quad M_j^2 = \frac{6e^{*j}}{d}, \quad j \in [h]_0, \quad \log h = \lceil \log(1 + \ln^*(d/6)) \rceil. \quad (31)$$

Therefore, by the previous lemma and Lemma 2.1, we get the following result.

Theorem 7.3. *For $r = dm$ with integer $m \geq h + \log h$, the protocol π_u^* with parameters as set in (31) satisfies*

$$\text{MSE}(\pi_u^*, \mathbf{x}, \mathbf{y}) \leq \sum_{i \in [n]} \frac{\Delta_i}{n} \cdot \frac{64\sqrt{3}}{n2^{r/(d(2+2\ln^*(d/6)))}},$$

for all \mathbf{x}, \mathbf{y} satisfying (2), for every $\Delta = (\Delta_1, \dots, \Delta_n)$.

Proof. Denote by \hat{x} the output of the protocol. Then, by Lemmas 2.1 and Lemma 7.2, we get

$$\begin{aligned} \mathbb{E} [\|\hat{x} - \bar{x}\|_2^2] &\leq \frac{1}{n^2} \sum_{i=1}^n \alpha(Q; \Delta_i) \\ &\leq \frac{16\sqrt{3}}{n^2 N} \sum_{i=1}^n \Delta_i, \end{aligned}$$

where the previous inequality is by Lemma 7.2. The proof is completed by using

$$N \geq \frac{2^{m/h}}{2^{1+(\log h)/h}} \geq \frac{2^{m/h}}{4} \geq \frac{2^{m/(2+2\ln^*(d/6))}}{4},$$

where the first inequality follows from using $\lfloor x \rfloor \geq x - 1$ for the floor function in the value of N in (31), the second follows from the fact that $\log x \leq x, \forall x \geq 0$, and the third follows from $\lceil x \rceil \leq x + 1$ for the ceil function in the value of h in (31). \square

8 Numerical Experiments

We empirically demonstrate the performance of our proposed quantizers on the following mean estimation task.

Each client i has a d -dimensional vector $x_i = \mu + U_i^c$, where μ in $[0, 1]^d$ is constant mean vector and U_i^c is a random vector whose each coordinate is a Uniform random variable in $[-\Delta'/2, \Delta'/2]$. The server has side information y_i corresponding to x_i , where $y_i = \mu + U_i^s$, and U^s , too, is a random vector whose each coordinate is a Uniform random variable in $[-\Delta'/2, \Delta'/2]$. Note that the distance between each coordinate of x_i and y_i is bounded by Δ' .

We compare three different mean estimation protocols. The first protocol is our first Wyner-Ziv estimator that uses RMQ for all the clients. Note that this protocol uses the knowledge of Δ' to set the values of RMQ. The second protocol is our universal Wyner-Ziv estimator which uses RDAQ for all the clients. Here, instead of vanilla RDAQ, we will use boosted RDAQ to make use of all the available precision. Recall that this particular protocol operates without the knowledge of Δ' . Our third protocol uses RATQ for all the clients, an efficient quantizer for the ℓ_2 ball [43]. Note that this protocol neither uses the side information y_i nor the distance between side information and the input vectors and will serve as a baseline. We evaluate the performance of our protocols by root mean square error (RMSE) between \bar{x} , the sample average of x_i s, and its estimate formed by the server \hat{x} .

We fix the number of clients $n = 10$. We conduct the experiments at dimensions $d = 512$ and $d = 1024$, and at two different precision levels: 6 bits per dimension and 10 bits per dimension. For all these four experiments we track the performance of our three quantization protocols by changing Δ' . All the experiments are averaged over ten runs for statistical consistency. Our implementation is available online at GitHub¹⁸.

We use the following parameters for all the quantizers. For RMQ, we set $\epsilon = \frac{2\Delta'}{31}$ and $\frac{2\Delta'}{511}$ for precision 6 bits and 10 bits, respectively. For RDAQ and RATQ, we first normalize the vectors $\{x_i, y_i\}_{i \in [n]}$ using an the bound $\sqrt{d}(1 + \Delta'/2)$ on their ℓ_2 -norm. Then, we set¹⁹ $h = 4$ to compute the different scales $M_{j,s}$ in (31) for dimensions $d = 512, 1024$. In addition, we choose $N = 1$ and $N = 3$ for implementing 6 bit and 10 bit Boosted RDAQ, respectively. The final estimate is obtained by multiplying back the decoded output with $\sqrt{d}(1 + \Delta'/2)$.

We see in Figures 1, 2, 3, and 4 that RMQ comfortably outperforms the other two quantizers at possible parameter choices. This is expected, since the RMSE of RMQ is directly proportional to Δ' , which is very small in our experiments. Another consequence of this relation to Δ' is that RMSE increases at a much faster rate with increase Δ' for RMQ than any other protocol. In other

¹⁸https://github.com/shubhamjha-46/WZ_estimators.

¹⁹For RATQ too, we set $h = 4$ (see [43] for more details).

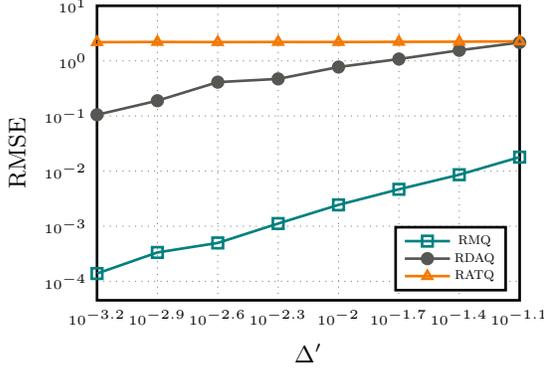


Figure 1: Comparison of RMQ, RDAQ, and RATQ at per coordinate precision of 6 bits and $d = 512$.

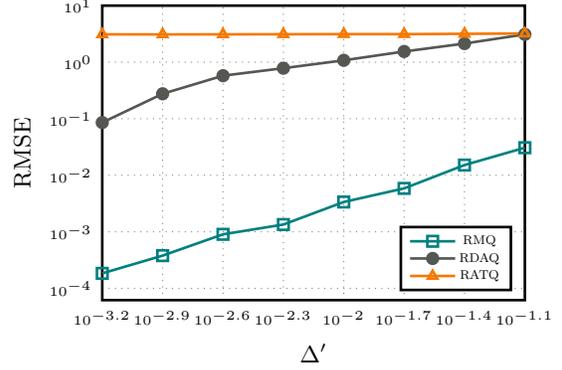


Figure 2: Comparison of RMQ, RDAQ, and RATQ at per coordinate precision of 6 bits and $d = 1024$.

words, the performance of RMQ will degrade at a much faster rate than RDAQ as the accuracy of side information degrades.

As can be seen in Figures 1, 2, RDAQ outperforms RATQ at 6 bits per dimension and both values of dimension. At precision level of 10 bits per dimension, however, RDAQ is better than RATQ at lower values of Δ' .

In other direction, we note that for all our protocols there is slight increase in RMSE for the same Δ' and bit precision, as the dimension increases from 512 to 1024. This is because ℓ_2 norm of the input and the ℓ_2 distance between input and side information depend on the dimension for our example, and our MSE upper bounds for all the quantizers depend on either one or both of these quantities.

Finally, we end with a remark on our choice of precision levels of 6 bits and 10 bits per dimension for this experiment. Notice that similar trends can be observed for precision levels lesser than dimension d . However, setting close to optimal parameters for these quantizers would have been much more tedious at precision levels lesser than the dimension. Since our experiment aimed to study the impact of side information on the accuracy of distributed mean estimation, we chose not to experiment with precision levels lesser than the dimension. The reason for not experimenting at 1 or 2 bits per dimension is that RDAQ is not operational below 6 bits per coordinate for the current dimension.

9 Proofs

9.1 Proof of Lemma 2.1

For the estimator \hat{x} in (4), with $\hat{x}_i = Q_i(x_i, y_i)$, we have

$$\mathbb{E} \left[\left\| \frac{1}{n} \cdot \sum_{i \in [n]} Q_i(x_i, y_i) - \frac{1}{n} \cdot \sum_{i \in [n]} x_i \right\|_2^2 \right]$$

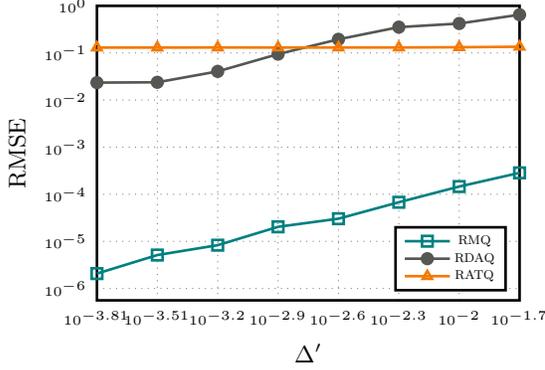


Figure 3: Comparison of RMQ, RDAQ, and RATQ at per coordinate precision of 10 bits and $d = 512$.

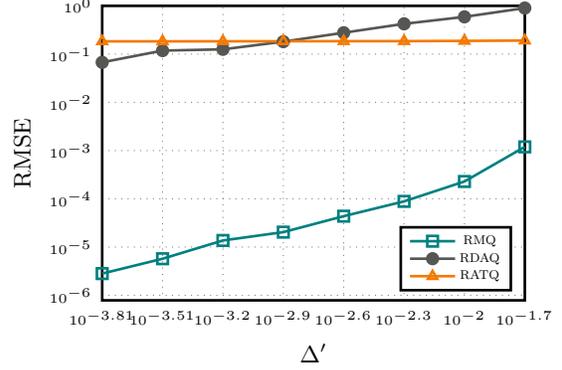


Figure 4: Comparison of RMQ, RDAQ, and RATQ at per coordinate precision of 10 bits and $d = 1024$.

$$\begin{aligned}
&= \frac{1}{n^2} \cdot \sum_{i \in [n]} \mathbb{E} [\|Q_i(x_i, y_i) - x_i\|_2^2] + \frac{1}{n^2} \cdot \sum_{i \neq j} \mathbb{E} [\langle Q_i(x_i, y_i) - x_i, Q_j(x_j, y_j) - x_j \rangle] \\
&= \frac{1}{n^2} \cdot \sum_{i \in [n]} \mathbb{E} [\|Q_i(x_i, y_i) - x_i\|_2^2] + \frac{1}{n^2} \cdot \sum_{i \neq j} \langle \mathbb{E} [Q_i(x_i, y_i)] - x_i, \mathbb{E} [Q_j(x_j, y_j)] - x_j \rangle \\
&= \frac{1}{n^2} \cdot \sum_{i \in [n]} \mathbb{E} [\|Q_i(x_i, y_i) - x_i\|_2^2] + \left(\frac{1}{n} \cdot \sum_i \|\mathbb{E} [Q_i(x_i, y_i)] - x_i\|_2 \right)^2 \\
&\quad - \frac{1}{n^2} \cdot \sum_i \|\mathbb{E} [Q_i(x_i, y_i)] - x_i\|_2^2 \\
&\leq \frac{1}{n^2} \cdot \sum_{i \in [n]} \mathbb{E} [\|Q_i(x_i, y_i) - x_i\|_2^2] + \frac{(n-1)}{n^2} \cdot \sum_i \|\mathbb{E} [Q_i(x_i, y_i)] - x_i\|_2^2,
\end{aligned}$$

where the second identity uses the independence of $Q_i(x_i, y_i)$ for different i and the final step uses Jensen's inequality. The result follows by bound each term using the fact that \mathbf{x} and \mathbf{y} satisfy (2) and the definitions of $\alpha(Q_i, \Delta_i)$ and $\beta(Q_i, \Delta_i)$, for $i \in [n]$. \square

9.2 Proof of Lemma 3.1

As mentioned in (5), the integer \tilde{z} found in Alg. 2 satisfies $\mathbb{E} [\tilde{z}\varepsilon] = x$ and $|x - \tilde{z}\varepsilon| < \varepsilon$. Therefore, it suffices to show that the output of the quantizer satisfies $Q_M(x, y) = \tilde{z}\varepsilon$.

To see that $Q_M(x, y) = \tilde{z}\varepsilon$, denote the lattice used in decoding Alg. 3 as $\mathbb{Z}_{w, \varepsilon} := \{(zk + w) \cdot \varepsilon : z \in \mathbb{Z}\}$. The decoding algorithm finds the point in $\mathbb{Z}_{w, \varepsilon}$ that is closest to y . Note that $w = \tilde{z} \bmod k$, whereby $\tilde{z}\varepsilon$ is a point in this lattice. Further, for any other point $\lambda \neq \tilde{z}\varepsilon$ in the lattice, we must have

$$|\lambda - \tilde{z}\varepsilon| \geq k\varepsilon,$$

and so, by triangular inequality, that

$$|\lambda - y| \geq |\lambda - \tilde{z}\varepsilon| - |\tilde{z}\varepsilon - y| \geq k\varepsilon - |\tilde{z}\varepsilon - y|.$$

Thus, $\tilde{z}\varepsilon$ is closer to y than λ if

$$k\varepsilon > 2|\tilde{z}\varepsilon - y|. \quad (32)$$

Next, by using (5) once again, we have

$$|\tilde{z}\varepsilon - y| \leq |\tilde{z}\varepsilon - x| + |x - y| < \varepsilon + \Delta',$$

which by condition (7) in the lemma implies that (32) holds. It follows that $|\lambda - y| > |\tilde{z}\varepsilon - y|$ for every $\lambda \in \mathbb{Z}_{w,\varepsilon}$, which shows that $Q_{\mathbf{M}}(x, y) = \tilde{z}\varepsilon$ and completes the proof. \square

9.3 Proof of Lemma 3.2

Recall from Remark 1 that for the random matrix R given in (8), for every vector $z \in \mathbb{R}^d$, the random variables $Rz(i)$, $i \in [d]$, are sub-Gaussian with variance parameter $\|z\|_2^2/d$. Furthermore, we need the following bound for “truncated moments” of sub-Gaussian random variables.

Lemma 9.1. *For a sub-Gaussian random Z with variance factor σ^2 and every $t \geq 0$, we have*

$$\mathbb{E} [Z^2 \mathbb{1}_{\{|Z|>t\}}] \leq 2(2\sigma^2 + t^2)e^{-t^2/2\sigma^2}.$$

Proof. Note that for any nonnegative random variable U , it can be verified that

$$\mathbb{E} [U \mathbb{1}_{\{U>x\}}] = xP(U > x) + \int_x^\infty P(U > u) du.$$

Upon substituting $U = Z^2$ and $x = t^2$, along with the fact that Z is sub-Gaussian with variance parameter σ^2 , we get

$$\begin{aligned} \mathbb{E} [Z^2 \mathbb{1}_{\{Z^2>t^2\}}] &= t^2 P(Z^2 > t^2) + \int_{t^2}^\infty P(Z^2 > u) du \\ &\leq 2t^2 e^{-t^2/2\sigma^2} + 2 \int_{t^2}^\infty e^{-u/2\sigma^2} du \\ &\leq 2(t^2 + 2\sigma^2)e^{-t^2/2\sigma^2}, \end{aligned}$$

which completes the proof. \square

We now handle the MSE $\alpha(Q)$ and bias $\beta(Q)$ separately below.

Bound for MSE $\alpha(Q)$: Denote by $Q_{\mathbf{M},R}(x, y)$ the final quantized value of the quantizer RMQ. For convenience, we abbreviate

$$\hat{x}_R := R Q_{\mathbf{M},R}(x, y).$$

Observe that $\hat{x}_R = \sum_{i \in [d]} Q_{\mathbf{M}}(Rx(i), Ry(i))e_i$, where $Q_{\mathbf{M}}$ is the MQ of Alg. 2 and 3 with parameters $k \geq$ and Δ' set as in the statement of the lemma. Since R is a unitary transform, we have

$$\mathbb{E} [\|Q_{\mathbf{M},R}(x, y) - x\|_2^2] = \mathbb{E} [\|\hat{x}_R - Rx\|_2^2]$$

$$\begin{aligned}
&= \sum_{i=1}^d \mathbb{E} [(\hat{x}_R(i) - Rx(i))^2] \\
&= \sum_{i=1}^d \mathbb{E} [(\hat{x}_R(i) - Rx(i))^2 \mathbb{1}_{\{|R(x-y)(i)| \leq \Delta'\}}] \\
&\quad + \sum_{i=1}^d \mathbb{E} [(\hat{x}_R(i) - Rx(i))^2 \mathbb{1}_{\{|R(x-y)(i)| \geq \Delta'\}}] \tag{33}
\end{aligned}$$

We consider each error term on the right-side above separately. We can view the first term as the error corresponding to MQ, when the input lies in its “acceptance range.” Specifically, under the event $\{|R(x-y)(i)| \leq \Delta'\}$, we get by Lemma 3.1 that

$$|\hat{x}_R(i) - Rx(i)| \leq \varepsilon = \frac{2\Delta'}{k-2}, \quad \text{almost surely,}$$

whereby

$$\sum_{i=1}^d \mathbb{E} [(\hat{x}_R(i) - Rx(i))^2 \mathbb{1}_{\{|R(x-y)(i)| \leq \Delta'\}}] \leq d\varepsilon^2. \tag{34}$$

The second term on the right-side of (33) corresponds to the error due to “overflow” and is handled using concentration bounds for the rotated vectors. Specifically, we get

$$\begin{aligned}
&\sum_{i=1}^d \mathbb{E} [(\hat{x}_R(i) - Rx(i))^2 \mathbb{1}_{\{|R(x-y)(i)| \geq \Delta'\}}] \\
&\leq 2 \sum_{i=1}^d [\mathbb{E} [(\hat{x}_R(i) - Ry(i))^2 \mathbb{1}_{\{|R(x-y)(i)| \geq \Delta'\}}] + \mathbb{E} [(Rx(i) - Ry(i))^2 \mathbb{1}_{\{|R(x-y)(i)| \geq \Delta'\}}]] \\
&\leq 2k^2\varepsilon^2 \sum_{i=1}^d P(|R(x-y)(i)| \geq \Delta') + 2 \sum_{i=1}^d \mathbb{E} [(Rx(i) - Ry(i))^2 \mathbb{1}_{\{|R(x-y)(i)| \geq \Delta'\}}] \\
&\leq 4dk^2\varepsilon^2 e^{-d\Delta'^2/2\Delta^2} + 2 \sum_{i=1}^d \mathbb{E} [(Rx(i) - Ry(i))^2 \mathbb{1}_{\{|R(x-y)(i)| \geq \Delta'\}}] \\
&\leq 4dk^2\varepsilon^2 e^{-d\Delta'^2/2\Delta^2} + 4(2\Delta^2 + d\Delta'^2) e^{-\frac{d\Delta'^2}{2\Delta^2}}, \tag{35}
\end{aligned}$$

where the second inequality follows upon noting that from the description decoder of MQ in Alg. 3 that $|\hat{x}_R(i) - Ry(i)| \leq \varepsilon k$ almost surely for each $i \in [d]$; the third inequality uses the fact that $R(x-y)(i)$ is sub-Gaussian with variance parameter $\|x-y\|_2^2/d \leq \Delta^2/d$; and fourth inequality is by Lemma 9.1.

Upon combining (33), (34), and (35), and substituting $\varepsilon = 2\Delta'/(k-2)$ and $\Delta'^2 = 6(\Delta^2/d) \log \Delta/\delta$, we obtain

$$\begin{aligned}
\mathbb{E} [\|Q_{M,R}(x,y) - x\|_2^2] &\leq d\varepsilon^2 + 4dk^2\varepsilon^2 e^{-\frac{d\Delta'^2}{2\Delta^2}} + 4(2\Delta^2 + d\Delta'^2) e^{-\frac{d\Delta'^2}{2\Delta^2}} \tag{36} \\
&= 24 \frac{\Delta^2}{(k-2)^2} \ln \frac{\Delta}{\delta} + 96\delta^2 \left(\frac{k}{k-2} \right)^2 \cdot \frac{\ln(\Delta/\delta)}{(\Delta/\delta)} + 8\delta^2 \cdot \frac{1+3\ln(\Delta/\delta)}{(\Delta/\delta)}
\end{aligned}$$

$$\leq 24 \frac{\Delta^2}{(k-2)^2} \ln \frac{\Delta}{\delta} + \left(\frac{96}{e} \left(\frac{k}{k-2} \right)^2 + \frac{24}{e^{2/3}} \right) \cdot \delta^2,$$

where we used $(1 + 3 \ln u)/u \leq 3/e^{2/3}$ and $(\ln u)/u \leq 1/e$ for every $u > 0$. We conclude by noting that for $k \geq 4$,

$$\left(\frac{96}{e} \left(\frac{k}{k-2} \right)^2 + \frac{24}{e^{2/3}} \right) \leq 154.$$

Bias $\beta(Q)$: The calculation for the bias is similar to that we used to bound the second term on the right-side of (33). Using the notation \hat{x}_R introduced above, we have

$$\begin{aligned} & \|\mathbb{E}[Q_{M,R}] - x\|_2 \\ &= \|\mathbb{E}[R^{-1}(\hat{x}_R - Rx)]\|_2 \\ &= \|R\mathbb{E}[R^{-1}(\hat{x}_R - Rx)]\|_2 \\ &= \|\mathbb{E}[RR^{-1}(\hat{x}_R - Rx)]\|_2 \\ &= \|\mathbb{E}[\hat{x}_R - Rx]\|_2, \end{aligned}$$

where the second identity holds since R is a unitary matrix.

Further, since $Q_M(x, y)$ is an unbiased estimate of x when $|x - y| \leq \Delta'$ (see Lemma 3.1), by (34) and (35) we obtain

$$\begin{aligned} \|\mathbb{E}[\hat{x}_R - Rx]\|_2^2 &\leq \sum_{i=1}^d \mathbb{E}[(\hat{x}_R(i) - Rx(i)) \mathbb{1}_{|R(x-y)_i| \geq \Delta'}]^2 \\ &\leq \sum_{i=1}^d \mathbb{E}[(\hat{x}_R(i) - Rx(i))^2 \mathbb{1}_{|R(x-y)(i)| \geq \Delta'}] \\ &\leq 154 \delta^2, \end{aligned}$$

which completes the proof. \square

9.4 Proof of Lemma 3.3

Mean Square Error $\alpha(Q_{S,R})$: From the description of Algorithms 6 and 7, we know that the quantized output of subsampled RMQ Q_{WZ} for an input x is

$$\begin{aligned} Q_{WZ}(x) &= R^{-1} \hat{x}_R, \text{ where} \\ \hat{x}_R &= \frac{1}{\mu} \sum_{i \in [d]} (Q_M(Rx(i), Ry(i)) - Ry(i)) \mathbb{1}_{\{i \in S\}} e_i + Ry, \end{aligned}$$

and $Q_M(Rx(i), Ry(i))$ denotes the quantized output of the modulo quantizer for an input $Rx(i)$ and side-information $Ry(i)$. Use the shorthand $Q(Rx(i))$ for $Q_M(Rx(i), Ry(i))$, we have

$$\mathbb{E}[\|Q_{WZ}(x) - x\|_2^2]$$

$$\begin{aligned}
&= \sum_{i \in [d]} \mathbb{E} \left[\left(\frac{1}{\mu} (Q(Rx(i)) - Ry(i)) \mathbb{1}_{\{i \in S\}} - (Rx(i) - Ry(i)) \right)^2 \right] \\
&\leq 2 \sum_{i \in [d]} \mathbb{E} \left[\frac{1}{\mu^2} (Q(Rx(i)) - Rx(i))^2 \mathbb{1}_{\{i \in S\}} \right] \\
&\quad + 2 \sum_{i \in [d]} \mathbb{E} \left[\left(\frac{1}{\mu} (Rx(i) - Ry(i)) \mathbb{1}_{\{i \in S\}} - (Rx(i) - Ry(i)) \right)^2 \right] \\
&= \sum_{i \in [d]} \frac{2}{\mu} \mathbb{E} \left[(Q(Rx(i)) - Rx(i))^2 \right] + 2 \sum_{i \in [d]} \mathbb{E} \left[(Rx(i) - Ry(i))^2 \right] \cdot \mathbb{E} \left[\left(\frac{1}{\mu} \mathbb{1}_{\{i \in S\}} - 1 \right)^2 \right] \\
&= \sum_{i \in [d]} \frac{2}{\mu} \mathbb{E} \left[(Q(Rx(i)) - Rx(i))^2 \right] + 2 \sum_{i \in [d]} \mathbb{E} \left[(Rx(i) - Ry(i))^2 \right] \cdot \frac{1 - \mu}{\mu} \\
&\leq \frac{2\alpha(Q_{M,R})}{\mu} + \frac{2\Delta^2}{\mu},
\end{aligned} \tag{37}$$

where we used the inequality: $(a + b)^2 \leq 2(a^2 + b^2)$, the independence of S and R in the second identity and used the fact that R is unitary in the final step.

Bias $\beta(Q_{S,R})$: This follows upon noting that the conditional expectation (over S) of the output of subsampled RMQ given R is the vector $R^{-1} \sum_{i \in [d]} Q_M(Rx(i), Ry(i))e_i$, which, in turn, is equivalent in distribution to the output of RMQ. \square

9.5 Proof of Theorem 3.5

We denote $\Delta_{min} = \min_{i \in [d]} \Delta_i$ and set y_i s to be 0. Let x_1, \dots, x_n be an *iid* sequence with common distribution such that for all $j \in [d]$ we have

$$x_1(j) = \begin{cases} \frac{\Delta_{min}}{\sqrt{d}} & \text{w.p. } \frac{1 + \alpha(j)\delta}{2} \\ -\frac{\Delta_{min}}{\sqrt{d}} & \text{w.p. } \frac{1 - \alpha(j)\delta}{2}, \end{cases}$$

where $\alpha \in \{-1, 1\}^d$ is generated uniformly at random. We have the following Lemma for such x_i s, which provides a lower bound for the MSE of any estimator of the mean of the distribution of x_i s.

Lemma 9.2. *For x_1, \dots, x_n generated as above and any estimator \hat{x} of the mean formed using only r -bit quantized version of x_i s, we have²⁰*

$$\mathbb{E} \left[\left\| \hat{x} - \frac{\delta \Delta_{min}}{\sqrt{d}} \alpha \right\|_2^2 \right] \geq c' \cdot \frac{d \Delta_{min}^2}{nr},$$

where $c' < 1$ is a universal constant.

Proof of Lemma 9.2 follows from either [16, Proposition 2] or [2, Theorem 11].

²⁰Note that the side information y_i s are all set to 0.

The proof of Theorem 3.5 is completed by using this claim. Specifically, using $2a^2 + 2b^2 \geq (a+b)^2$, we have

$$2\mathbb{E} [\|\hat{x} - \bar{x}\|_2^2] + 2\mathbb{E} \left[\left\| \bar{x} - \frac{\delta\Delta_{min}}{\sqrt{d}}\alpha \right\|_2^2 \right] \geq \mathbb{E} \left[\left\| \hat{x} - \frac{\delta\Delta_{min}}{\sqrt{d}}\alpha \right\|_2^2 \right],$$

which, along with the observation that

$$\mathbb{E} \left[\left\| \bar{x} - \frac{\delta\Delta_{min}}{\sqrt{d}}\alpha \right\|_2^2 \right] \leq \frac{\Delta_{min}^2}{n},$$

gives

$$\begin{aligned} \mathbb{E} [\|\hat{x} - \bar{x}\|_2^2] &\geq \frac{c'd\Delta_{min}^2}{2nr} - \frac{\Delta_{min}^2}{n} \\ &\geq \frac{c'\Delta_{min}^2 d}{4nr}, \end{aligned}$$

when $(d/r) \geq 4/c'$. The proof is completed by setting $c = c'/4$. \square

Remark 9. Since the lower bound in [2] holds for sequentially interactive protocols, if we allow interactive protocols for mean estimation where client i gets to see the messages transmitted by the clients j in $[i-1]$, and can design its quantizers based on these previous messages, even then the lower bound above will hold.

9.6 Proof of Lemma 4.1

We will prove a general result which will not only prove Lemma 4.1 but will also be useful in the proof of Lemma 4.2. Consider x and y in \mathbb{R}^d such that each coordinate of both x and y lies in $[-M, M]$. Also, consider the following generalization of DAQ:

$$Q_D(x, y) = \sum_{i=1}^d 2M (\mathbb{1}_{\{U_i \leq x(i)\}} - \mathbb{1}_{\{U_i \leq y(i)\}}) e_i + y,$$

where $\{U_i\}_{i \in [d]}$ are iid uniform random variables in $[-M, M]$. We will show that

$$\mathbb{E} [Q_D(x, y)] = x \quad \text{and} \quad \mathbb{E} [\|Q_D(x, y) - x\|_2^2] \leq 2M\|x - y\|_1, \quad (38)$$

which upon setting $M = 1$ proves Lemma 4.1.

Towards proving (38), note that from the estimate formed by Q_D , it is easy to see that $\mathbb{E} [Q_D(x, y)] = x$. The MSE can be bounded as follows:

$$\begin{aligned} \mathbb{E} [\|Q_D(x, y) - x\|_2^2] &= \sum_{i=1}^d \mathbb{E} [(2M (\mathbb{1}_{\{U_i \leq x(i)\}} - \mathbb{1}_{\{U_i \leq y(i)\}}) - (x(i) - y(i)))^2] \\ &= \sum_{i=1}^d 4M^2 \frac{|x(i) - y(i)|}{2M} - \|x - y\|_2^2 \\ &= 2M\|x - y\|_1 - \|x - y\|_2^2, \end{aligned}$$

where we used the observations that $2M (\mathbb{1}_{\{U_i \leq x(i)\}} - \mathbb{1}_{\{U_i \leq y(i)\}})$ is an unbiased estimate of $(x(i) - y(i))$ and that $(\mathbb{1}_{\{U_i \leq x(i)\}} - \mathbb{1}_{\{U_i \leq y(i)\}})^2$ equals one if and only if exactly one of the indicators is one, which in turn happens with probability $\frac{|x(i) - y(i)|}{2M}$. \square

9.7 Proof of Lemma 4.2

Worst-case bias $\beta(Q_{\mathbf{d},R}\Delta)$: Since the final interval $[-M_{h-1}, M_{h-1}]$ contains $[-1, 1]$, we can see that $\mathbb{E}[Q_{\mathbf{d},R}(x, y)] = x$.

Worst-case MSE $\alpha(Q_{\mathbf{d},R}; \Delta)$: We denote by B_{ij}^x and B_{ij}^y the bits

$$B_{ij}^x = \mathbb{1}_{\{U(i,j) \leq Rx(i)\}} \quad \text{and} \quad B_{ij}^y = \mathbb{1}_{\{U(i,j) \leq Ry(i)\}}.$$

Then, the final quantized value of the quantizer RDAQ can be expressed as $Q_{\mathbf{d},R}(X) = R^{-1}\hat{x}_R$ where, with $z^*(i)$ denoting the smallest M_j such that the interval $[-M_j, M_j]$ contains $Rx(i)$ and $Ry(i)$ and $[h]_0 = \{0, \dots, h-1\}$,

$$\hat{x}_R := \sum_{i \in \{1, \dots, d\}} \left(\sum_{j \in [h]_0} 2M_j \cdot (B_{ij}^x - B_{ij}^y) + Ry(i) \right) \mathbb{1}_{\{z^*(i)=j\}} e_i.$$

Since R is a unitary transform, we get

$$\begin{aligned} \mathbb{E} [\|Q_{\mathbf{d},R}(x) - x\|_2^2] &= \mathbb{E} [\|RQ_{\mathbf{d},R}(x) - Rx\|_2^2] \\ &= \mathbb{E} [\|\hat{x}_R - Rx\|_2^2] \\ &= \sum_{i \in [d]} \mathbb{E} [(\hat{x}_R(i) - Rx(i))^2] \\ &= \sum_{i \in [d]} \mathbb{E} \left[\left(\sum_{j \in [h]_0} (2M_j \cdot (B_{ij}^x - B_{ij}^y) + Ry(i) - Rx(i)) \mathbb{1}_{\{z^*(i)=j\}} \right)^2 \right] \\ &= \sum_{i \in [d]} \sum_{j \in [h]_0} \mathbb{E} \left[(2M_j (B_{ij}^x - B_{ij}^y) + Ry(i) - Rx(i))^2 \mathbb{1}_{\{z^*(i)=j\}}, \right] \end{aligned}$$

where the last identity uses $\mathbb{1}_{\{z^*(i)=j_1\}} \mathbb{1}_{\{z^*(i)=j_2\}} = 0$ for all $j_1 \neq j_2$, to cancel the cross-terms in the expansion of $(\hat{x}_R(i) - Rx(i))^2$. Conditioning on R and using the independence of $\mathbb{1}_{\{z^*(i)=j\}}$ from the randomness used in MQ, we get

$$\begin{aligned} \mathbb{E} [\|Q_{\mathbf{d},R}(x) - x\|_2^2] &= \sum_{i \in [d]} \sum_{j \in [h]_0} \mathbb{E} \left[\mathbb{E} \left[(2M_j (B_{ij}^x - B_{ij}^y) + Ry(i) - Rx(i))^2 \mid R \right] \mathbb{1}_{\{z^*(i)=j\}} \right] \\ &\leq \sum_{i \in [d]} \sum_{j \in [h]_0} \mathbb{E} [2M_j |Rx(i) - Ry(i)| \mathbb{1}_{\{z^*(i)=j\}}], \\ &\leq \sum_{i \in [d]} \mathbb{E} [2M_0 |Rx(i) - Ry(i)| \mathbb{1}_{\{z^*(i)=0\}}] \\ &\quad + \sum_{i \in [d]} \sum_{j \in [h-1]} \mathbb{E} [2M_j |Rx(i) - Ry(i)| \mathbb{1}_{\{z^*(i)=j\}}], \end{aligned}$$

$$\begin{aligned}
&\leq \sum_{i \in [d]} \mathbb{E} [2M_0 |Rx(i) - Ry(i)|] \\
&\quad + \sum_{i \in [d]} \sum_{j \in [h-1]} \mathbb{E} [2M_j |Rx(i) - Ry(i)| \mathbb{1}_{\{z^*(i)=j\}}], \tag{39}
\end{aligned}$$

where the first inequality follows from (38) in the proof of Lemma 4.1. Next, noting that

$$\mathbb{1}_{\{z^*(i)=j\}} \leq \mathbb{1}_{\{|RX(i)| \geq M_{j-1}\}} + \mathbb{1}_{\{|RY(i)| \geq M_{j-1}\}} \quad \text{almost surely,}$$

an application of the Cauchy-Schwarz inequality yields

$$\begin{aligned}
&\mathbb{E} [2M_j |Rx(i) - Ry(i)| \mathbb{1}_{\{z^*(i)=j\}}] \\
&\leq 2M_j \mathbb{E} [(Rx(i) - Ry(i))^2]^{1/2} \mathbb{E} [(\mathbb{1}_{\{|RX(i)| \geq M_{j-1}\}} + \mathbb{1}_{\{|RY(i)| \geq M_{j-1}\}})^2]^{1/2} \\
&\leq 2M_j \mathbb{E} [(Rx(i) - Ry(i))^2]^{1/2} (2P(|Rx(i)| \geq M_{j-1}) + 2P(|Ry(i)| \geq M_{j-1}))^{1/2} \\
&\leq 2M_j \mathbb{E} [(Rx(i) - Ry(i))^2]^{1/2} \left(8e^{-\frac{dM_{j-1}^2}{2}}\right)^{1/2}, \tag{40}
\end{aligned}$$

where the second inequality uses $(a+b)^2 \leq 2a^2 + 2b^2$ and the third uses subgaussianity of $Rx(i)$ and $Ry(i)$.

Substituting the upper bound in (40) for the second term in the RHS of (39) and using $\mathbb{E}[X] \leq \mathbb{E}[X^2]^{1/2}$ for the first term, we get

$$\begin{aligned}
\mathbb{E} [\|Q_{D,R}(x) - x\|_2^2] &\leq \sum_{i \in [d]} \mathbb{E} [|Rx(i) - Ry(i)|^2]^{1/2} \left(2M_0 + \sum_{j \in [h-1]} 2M_j \cdot \left(8e^{-\frac{dM_{j-1}^2}{2}}\right)^{1/2}\right) \\
&\leq \sqrt{d \cdot \mathbb{E} [\|Rx - Ry\|_2^2]} \left(2M_0 + \sum_{j \in [h-1]} 2M_j \cdot \left(8e^{-\frac{dM_{j-1}^2}{2}}\right)^{1/2}\right) \\
&= \sqrt{d \cdot \|x - y\|_2^2} \left(2M_0 + \sum_{j \in [h-1]} 2M_j \cdot \left(8e^{-\frac{dM_{j-1}^2}{2}}\right)^{1/2}\right) \\
&= \sqrt{d \cdot \|x - y\|_2^2} \left(2\sqrt{\frac{6}{d}} + \sum_{j \in [h-1]} 2\sqrt{\frac{6e^{*j}}{d}} \cdot \left(8e^{-1.5e^{*(j-1)}}\right)\right) \\
&= 8\sqrt{3} \cdot \sqrt{\|x - y\|_2^2} \left(1 + \sum_{j \in [h-1]} e^{-0.5e^{*(j-1)}}\right) \\
&\leq 16\sqrt{3} \cdot \sqrt{\|x - y\|_2^2},
\end{aligned}$$

where the second inequality uses the fact that $\sum_i \|a\|_1 \leq \sqrt{d}\|a\|_2$, the first and second identities follow from the fact that R is unitary transform and substituting for M_i s, the final inequality follows

from the bound of 1 for $\sum_{j=1}^{\infty} e^{-0.5e^{*(j-1)}}$, which, in turn, can be seen as follows

$$\begin{aligned}
e^{-0.5e^{*(j-1)}} &= e^{-0.5} + e^{-0.5e} + e^{-0.5e^e} + \sum_{j=3}^{\infty} e^{-0.5e^{*(j)}} \\
&\leq e^{-0.5} + e^{-0.5e} + e^{-0.5e^e} + \sum_{j=3}^{\infty} e^{-0.5je^e} \\
&\leq e^{-0.5} + e^{-0.5e} + e^{-0.5e^e} + \frac{1}{e^{e^e} - 1} \\
&\leq 1.
\end{aligned}$$

□

9.8 Proof of Lemma 4.3

Worst-case bias $\beta(Q_{\text{wz},u}; \Delta)$: It is straightforward to see that $\mathbb{E}[Q_{\text{wz},u}(x)] = x$.

Worst-case MSE $\alpha(Q_{\text{wz},u}; \Delta)$: We denote by B_{ij}^x and B_{ij}^y the bits

$$B_{ij}^x = \mathbb{1}_{\{U(i,j) \leq Rx(i)\}} \quad \text{and} \quad B_{ij}^y = \mathbb{1}_{\{U(i,j) \leq Ry(i)\}}.$$

Then, the quantized output can be stated as follows: noting that $Q_{\text{wz},u}(x) = R^{-1}\hat{x}_R$ where, with $z^*(i)$ denoting the smallest M_j such that the interval $[-M_j, M_j]$ contains $Rx(i)$ and $Ry(i)$,

$$\hat{x}_R := \left(\sum_{i \in \{1, \dots, d\}} \sum_{j \in \{0, \dots, h-1\}} 2M_j \cdot (B_{ij}^x - B_{ij}^y) \mathbb{1}_{\{z^*(i)=j\}} \mathbb{1}_{\{i \in S\}} \cdot e_i + Ry \right),$$

Since R is a unitary transform, the mean square error between $Q_{\text{wz},u}(x)$ and x can be bounded as in the proof of Lemma 4.2 as follows:

$$\begin{aligned}
\mathbb{E} [\|Q_{\text{wz},u}(x) - x\|_2^2] &= \mathbb{E} [\|\hat{x}_R - Rx\|_2^2] \\
&= \mathbb{E} [\|\hat{x}_R - Rx\|_2^2] \\
&= \sum_{i \in [d]} \mathbb{E} [\hat{x}_R(i) - Rx(i)]^2 \\
&= \sum_{i \in [d]} \sum_{j \in [h]} \mathbb{E} \left[(2M_j (B_{ij}^x - B_{ij}^y) \mathbb{1}_{\{i \in S\}} + Ry(i) - Rx(i))^2 \mathbb{1}_{\{z^*(i)=j\}} \right] \\
&= \sum_{i \in [d]} \sum_{j \in [h]} \mathbb{E} \left[\mathbb{E} \left[(2M_j (B_{ij}^x - B_{ij}^y) \mathbb{1}_{\{i \in S\}} + Ry(i) - Rx(i))^2 \mid R \right] \mathbb{1}_{\{z^*(i)=j\}} \right] \\
&\leq \sum_{i \in [d]} \sum_{j \in [h]} \mathbb{E} \left[\frac{2M_j}{\mu} \cdot |Rx(i) - Ry(i)| \cdot \mathbb{1}_{\{z^*(i)=j\}} \right],
\end{aligned}$$

where the inequality follows from similar calculations in the proof of Lemma 4.1. The rest of the analysis proceeds as that in the proof of Lemma 4.2. □

9.9 Proof of Theorem 5.1

Note that affine functions are 0-smooth and admitted in the class of L smooth functions. We use affine functions as difficult functions and follow the general recipe of [3], which in turn builds on [1, Section 4.5] and [2], to show the lower bounds for convex, Lipschitz optimization under communication constraints. The difficult functions we construct are the same as in many existing lower bounds for convex functions such as [5]. We consider the domain $\mathcal{X} = \{x \in \mathbb{R}^d : \|x\|_\infty \leq D/(2\sqrt{d})\}$, and consider the following class of functions on \mathcal{X} : For $v \in \{-1, 1\}^d$, let

$$f_v(x) := \frac{2\sigma\delta}{\sqrt{d}} \sum_{i=1}^d \left| x(i) - \frac{v(i)D}{2\sqrt{d}} \right|, \quad \forall x \in \mathcal{X}, \quad (41)$$

and x_v^* be its minimizer. Note that the gradient $g_v(x)$ of f_v at $x \in \mathcal{X}$ is equal to $-2\sigma\delta v/\sqrt{d}$, i.e., constant $\forall x$. For each f_v in (41), consider a sequence of n clients \mathbf{C} that output d -dimensional gradient vectors $\{\hat{g}_i(x_t)\}_{i \in [n]}$, each of whose coordinates takes value $-\sigma/\sqrt{d}$ or σ/\sqrt{d} independently with probabilities $(1 + 2\delta v(i))/2$ and $(1 - 2\delta v(i))/2$, respectively. The parameter $\delta > 0$ is to be chosen later. Note that the above client construction satisfies the set of assumptions in (18), (19) and (20).

Draw $V \sim \text{Unif}\{-1, 1\}^d$. With respect to the associated random function f_V , each client $\mathbf{C}1_i$ chooses a quantizer $Q_{i,t}$ to generate output $Q_{i,t}(\hat{g}_i(x_t))$. Denote by $Q^{nT} = (\{Q_{i,1}, \dots, Q_{i,T}\}_{i \in [n]})$ the vector of quantized outputs observed at the server. The following lower bound can be established by using results from [3, Lemma 3, 4]:

$$\mathbb{E}[f_V(\bar{x}_T) - f_V(x_V^*)] \geq \frac{D\sigma\delta}{3} \left[1 - \sqrt{\frac{2}{d} \sum_{j=1}^d I(V(j) \wedge Q^{nT})} \right]. \quad (42)$$

It remains to bound the mutual-information term for which one can use the independence across the clients and derive the following data-processing inequality based on the other techniques from [3]:

$$\sum_{j=1}^d I(V(j) \wedge Q^{nT}) \leq 29nT\delta^2(d \wedge r),$$

where $\delta \in (0, 1/6)$. Combining this with (42) and setting $\delta = \sqrt{d/(232(d \wedge r)nT)}$, we finally get

$$\mathbb{E}[f_V(\bar{x}_T) - f_V(x_V^*)] \geq \frac{1}{12\sqrt{58}} \frac{D\sigma}{\sqrt{nT}} \sqrt{\frac{d}{d \wedge r}},$$

where we need $T \geq d/(6nr)$ in order to enforce $\delta \leq 1/6$. The proof is completed by noting that $\mathcal{E}^*(\mathcal{X}, \mathcal{O}_{\text{sc}}, T, \mathcal{Q}_r) \geq \mathbb{E}[f_V(\bar{x}_T) - f_V(x_V^*)]$.

9.10 Proof of Lemma 5.2

Define $x^* = \text{argmin}_{x \in \mathcal{X}} f(x)$. We have that

$$\mathbb{E}[f(x_{t+1}) - x^*] = \mathbb{E}[f(x_{t+1}) - f(x_t)] + \mathbb{E}[f(x_t) - x^*]. \quad (43)$$

By smoothness,

$$\begin{aligned}
\mathbb{E}[f(x_{t+1}) - f(x_t)|x_t] &\leq \nabla f(x_t)^\top \mathbb{E}[x_{t+1} - x_t|x_t] + \frac{L}{2} \mathbb{E}[\|x_{t+1} - x_t\|^2|x_t] \\
&= -\eta \nabla f(x_t)^\top \mathbb{E}[\mathcal{M}(C_t^n)|x_t] + \frac{L\eta^2}{2} \mathbb{E}[\|\mathcal{M}(C_t^n)\|^2|x_t] \\
&\leq -\eta \nabla f(x_t)^\top \mathbb{E}[\mathcal{M}(C_t^n)|x_t] + \frac{\eta}{2} \mathbb{E}[\|\mathcal{M}(C_t^n)\|^2|x_t] \\
&= -\frac{\eta}{2} \|\nabla f(x_t)\|^2 + \frac{\eta}{2} \mathbb{E}[\|\mathcal{M}(C_t^n) - \nabla f(x_t)\|^2|x_t],
\end{aligned}$$

which further using the definition of α' in (21) and the law of total expectation imply

$$\mathbb{E}[f(x_{t+1}) - f(x_t)] \leq -\frac{\eta}{2} \mathbb{E}[\|\nabla f(x_t)\|^2] + \frac{\eta}{2} \alpha'^2(\mathcal{M}). \quad (44)$$

By convexity,

$$\begin{aligned}
\mathbb{E}[f(x_t) - x^*] &\leq \mathbb{E}[\nabla f(x_t)^\top (x_t - x^*)] \\
&= \mathbb{E}[(\nabla f(x_t) - \mathcal{M}(C_t^n))^\top (x_t - x^*)] + \mathbb{E}[\mathcal{M}(C_t^n)^\top (x_t - x^*)] \\
&= \mathbb{E}[(\nabla f(x_t) - \mathcal{M}(C_t^n))^\top (x_t - x^*)] + \frac{1}{2\eta} \mathbb{E}[\eta^2 \|\mathcal{M}(C_t^n)\|^2] \\
&\quad + \mathbb{E}[\|x_t - x^*\|^2 - \|x_t - \eta \mathcal{M}(C_t^n) - x^*\|^2] \\
&\leq \mathbb{E}[(\nabla f(x_t) - \mathcal{M}(C_t^n))^\top (x_t - x^*)] + \frac{1}{2\eta} \mathbb{E}[\eta^2 \|\mathcal{M}(C_t^n)\|^2] \\
&\quad + \mathbb{E}[\|x_t - x^*\|^2 - \|\Gamma_{\mathcal{X}}(x_t - \eta \mathcal{M}(C_t^n)) - x^*\|^2] \\
&\leq \beta'(\mathcal{M}) \cdot D + \frac{\eta}{2} \alpha'^2(\mathcal{M}) + \frac{\eta}{2} \mathbb{E}[\|\nabla f(x_t)\|^2] + \eta B \cdot \beta'(\mathcal{M}) \\
&\quad + \frac{1}{2\eta} \mathbb{E}[\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2],
\end{aligned} \quad (45)$$

where second inequality is due to a well known property of the projection operator $\Gamma_{\mathcal{X}}$ (see, for instance, Lemma 3.1, [13]), third inequality follows from Cauchy-Schwarz inequality and using the definitions in (21) and (22). Plugging (48) and (49) in (43), we have

$$\mathbb{E}[f(x_{t+1}) - x^*] \leq \beta'(\mathcal{M}) \cdot (D + \eta B) + \eta \cdot \alpha'^2(\mathcal{M}) + \frac{1}{2\eta} \mathbb{E}[\|x_t - x^*\|^2 - \|x_{t+1} - x^*\|^2].$$

Summing from $t = 0$ to $T - 1$, dividing by T , using the assumption that the domain \mathcal{X} has diameter at most D , and setting η as provided, the proof is completed. This general convergence bound will be used in our upper bound proofs below.

9.11 Proof of Theorem 5.3

From [43, Theorem 3.7], we use the following result.

Lemma 9.3. *Let Q_{RATQ} be the subsampled version of RATQ using $r \geq 3 + \lceil \log(1 + \ln^*(d/3)) \rceil$ bits. Then for Y such that $\|Y\|_2 \leq B^2$, we have*

$$\mathbb{E}[Q_{\text{RATQ}}(Y) | Y] = Y \quad \text{and} \quad \mathbb{E}[\|Q_{\text{RATQ}}(Y) - Y\|^2] \leq \frac{dB^2}{\frac{r}{3 + \lceil \log(1 + \ln^*(d/3)) \rceil} - 1}.$$

Further, for $t \in [T]$, we have $\mathcal{M}(C_t^n) = \frac{1}{n} \sum_{i \in [n]} Q_{\text{RATQ}}(\hat{g}_i(x_t))$ as in (23). Thus we have,

$$\alpha'^2(\mathcal{M}) \leq \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i \in [n]} Q_{\text{RATQ}}(\hat{g}_i(x_t)) - \frac{1}{n} \sum_{i \in [n]} \hat{g}_i(x_t) \right\|^2 \right] + \mathbb{E} \left[\left\| \frac{1}{n} \sum_{i \in [n]} \hat{g}_i(x_t) - \nabla f(x_t) \right\|^2 \right].$$

Since $\mathcal{M}(C_t^n)$ is an unbiased estimate, $\beta'(\mathcal{M}) = 0$. The proof is completed by bounding the two terms in the right-side above followed by using Lemma 5.2, which we do as follows. From Lemma 9.3, it follows that for $t \in [T]$, we have

$$\mathbb{E} \left[\left\| \frac{1}{n} \sum_{i \in [n]} Q_{\text{RATQ}}(\hat{g}_i(x_t)) - \frac{1}{n} \sum_{i \in [n]} \hat{g}_i(x_t) \right\|^2 \right] \leq \frac{dB^2}{n \left(\frac{r}{3 + \lceil \log(1 + \ln^*(d/3)) \rceil} - 1 \right)}.$$

From (27), we have

$$\mathbb{E} \left[\left\| \frac{1}{n} \sum_{i \in [n]} \hat{g}_i(x_t) - \nabla f(x_t) \right\|^2 \right] \leq \frac{\sigma^2}{n}.$$

9.12 Proof of Theorem 5.4

Subgaussian and subexponential norms. For our analysis, it will be convenient to recall the definition of subgaussian²¹ and subexponential norms of a random variable.

Definition 9.4 ([56]). A subgaussian norm of a subgaussian random variable X , denoted $\|X\|_{\psi_2}$, is defined as $\|X\|_{\psi_2} := \inf\{t > 0 : \mathbb{E} [e^{X^2/t^2}] \leq 2\}$. It follows that for a centered subgaussian random variable X , $\Pr(|X| \geq t) \leq 2e^{-\frac{t^2}{\|X\|_{\psi_2}^2}}$.

Definition 9.5 ([56, Def. 2.7.5]). A subexponential norm of a subexponential random variable X , denoted $\|X\|_{\psi_1}$, is defined as $\|X\|_{\psi_1} := \inf\{t > 0 : \mathbb{E} [e^{\frac{|X|}{t}}] \leq 2\}$. It follows that for a centered subexponential random variable X , $\Pr(|X| \geq t) \leq 2e^{-\frac{t}{\|X\|_{\psi_1}}}$.

Side information is close to gradient estimates. We begin by noting that side-information Y^{22} is close to the stochastic gradient estimates computed by clients in \mathcal{C}_2 . Specifically, setting the parameters as $\log \ell_1 = \lceil \log \frac{2B}{\sigma} + 1 \rceil$ and $r_1 = r / \lceil \log \frac{2B}{\sigma} + 1 \rceil$ for clients in \mathcal{C}_1 , we get the following.

Lemma 9.6. For all $x \in \mathbb{R}^d$, $j \in \mathcal{C}_2$, $i \in [d]$, and a universal constant $c_3 > 0$, we have

$$\Pr(|R\hat{g}_j(x)(i) - RY(i)| \geq t) \leq 2e^{-c_3 \min\{\frac{t^2}{\sigma^2}, \frac{t\sqrt{d}}{\sigma}\}} + 2e^{-c_3 \frac{t^2 d}{\sigma^2}},$$

where R is a random Hadamard matrix (8) and for another universal constant $c_4 > 0$,

$$\sigma'^2 = \frac{c_4 8d\sigma^2 \lceil \log(2B/\sigma + 1) \rceil}{nr}. \quad (46)$$

²¹ $\|\cdot\|_{\psi_2}$ is indeed a norm.

²²For convenience, we drop the iteration subscript t in this Section.

Remark 10. In the analysis for RMQ presented in Section 3.2, the difference between the coordinates of the rotated input and rotated side information had subgaussian tails. However, note that in Lemma 9.6, we can only prove a slightly weaker concentration result.

Towards proving Lemma 9.6, we begin by showing the following result which holds from the subgaussian properties of uniform quantizer error and standard properties of subgaussian random variables.

Lemma 9.7. *For all $x \in \mathbb{R}^d$ and $i \in [d]$ we have*

$$\|Y(i) - \nabla f(x)(i)\|_{\psi_2}^2 \leq \sigma'^2.$$

Proof. We will prove the theorem for $Y(1)$ since the argument remains the same for all $Y(i)$ s. From the description of CUQ, we note that $Q_u(\hat{g}_j(x)(1))$ satisfies

$$\|Q_u(\hat{g}_j(x)(1)) - \hat{g}_j(x)(1)\|_{\psi_2}^2 \leq \frac{4c_4B^2}{(\ell_1 - 1)^2}, \quad \forall j \in S_1,$$

for some universal constant $c_4 > 0$. Also, from (19), we have $\|\hat{g}_j(x)(1) - \nabla f(x)(1)\|_{\psi_2}^2 \leq c_4\sigma^2$ for the same constant c_4 above. Further, using the triangle inequality and the fact that $(a+b)^2 \leq 2a^2 + 2b^2$,

$$\|Q_u(\hat{g}_j(x)(1)) - \nabla f(x)(i)\|_{\psi_2}^2 \leq \frac{c_48B^2}{(\ell_1 - 1)^2} + 2c_4\sigma^2.$$

The proof is completed upon noting that the average of N *iid* zero mean subgaussian random variables $\{X_i\}_{i \in [N]}$ has a subgaussian norm square equal to $\|X_1\|_{\psi_2}^2/N$ and the fact that we use $N = nr_1/(2d)$ samples to form $Y(1)$. \square

Remark 11. In order to quantize a d -dimensional gradient to $r \leq d$ bits, the technique of uniform sampling has been used in recent papers on distributed optimization (*cf.* [55], [43]). However, notice that these works merely required the quantized gradient estimate to be close to the true gradient in mean square sense. In our case, in order to leverage our Wyner-Ziv compression algorithms, we need side-information to be close to the true gradient in a much stronger sense. Therefore, we refrain from using uniform sampling and instead use the clients to quantize separate, smaller blocks of coordinates.

Rotation of side-information is close to the rotation of true gradient. Using standard properties of subgaussian random variables (see [56, Lemma 2.7.7 and Theorem 2.8.1]), we can show the following.

Lemma 9.8. *For all $x \in \mathbb{R}^d$ and $i \in [d]$ we have for a universal constant $c_5 > 0$*

$$\Pr(|RY(i) - R\nabla f(x)(i)| \geq t) \leq 2e^{(-c_5 \min\{t^2/\sigma'^2, t\sqrt{d}/\sigma'\})}.$$

Proof. The proof follows from combining two facts. First, note that for a sequence $\{X_i\}_{i \in [N]}$ of zero mean, *iid* subexponential random variables, we have from [56, Theorem 2.8.1]

$$\Pr\left\{\left|\sum_{i=1}^N X_i\right| \geq t\right\} \leq 2e^{-c_5 \min\left\{\frac{t^2}{N\|X_1\|_{\psi_1}^2}, \frac{t}{\|X_1\|_{\psi_1}}\right\}},$$

for some universal constant $c_5 > 0$.

Also, note that the product of two subgaussian random variables X and Y is subexponential random variable with subexponential norm bounded as follows (see [56, Lemma 2.7.7]): $\|XY\|_{\psi_1} \leq \|X\|_{\psi_2}\|Y\|_{\psi_2}$. Notice that $|RY(i) - R\nabla f(x)(i)| = \sum_{i \in [d]} U(i)V(i)$, where $U(i), V(i)$ are zero mean, *iid*, subgaussian random variables with subgaussian norms as $1/\sqrt{d}$ and σ' (*cf.* Lemma 9.7). \square

Finally, proceeding in the same manner as in [43, Lemma 5.8], we can show that the coordinates of the rotated stochastic gradient are close to the coordinates of the rotated true gradient.

Lemma 9.9. *For all $x \in \mathbb{R}^d$ and the universal constant $c_4 > 0$ as in Lemma 9.7, we have*

$$\|R\hat{g}_j(x)(i) - R\nabla f(x)(i)\|_{\psi_2}^2 \leq c_1\sigma^2/d.$$

Thus, random rotation allows us to convert the ℓ_2 norm bound in assumption (19) to a ℓ_∞ bound.

We now choose $c_3 = \min\{c_4, c_5\}$. Using the inequality: $\max(a, b) \leq a + b$ and the property of subgaussian random variable in Definition 9.4, Lemma 9.6 follows from combining Lemmas 9.8 and 9.9. Next, we present a lemma similar to Lemma 9.1 towards evaluating bounds on $\alpha'(\mathcal{M})$ and $\beta'(\mathcal{M})$.

Lemma 9.10. *For a random variable Z such that*

$$\Pr(|Z| \geq t) \leq 2e^{-\frac{c_3 t^2}{\sigma'^2}} + 2e^{-\frac{c_3 t \sqrt{d}}{\sigma'}} + 2e^{-\frac{c_3 t^2 d}{\sigma^2}},$$

where $c_3 > 0$ is some universal constant, we have

$$\mathbb{E}[Z^2 \mathbb{1}_{\{|Z| > t\}}] \leq 2 \left(\frac{\sigma'^2}{c_3} + t^2 \right) e^{-\frac{c_3 t^2}{\sigma'^2}} + 2 \left(\frac{\sigma^2}{dc_3} + t^2 \right) e^{-\frac{c_3 t^2 d}{\sigma^2}} + 2 \left(\frac{2\sigma'^2}{dc_3^2} + \frac{2\sigma' t}{c_3 \sqrt{d}} + t^2 \right) e^{-\frac{c_3 t \sqrt{d}}{\sigma'}}.$$

Proof. For any nonnegative random variable U , it can be seen

$$\mathbb{E}[U \mathbb{1}_{\{U > x\}}] = x \Pr(U > x) + \int_x^\infty \Pr(U > u) du.$$

Upon substituting $U = Z^2$ and $x = t^2$, along with the fact that Z has the tail behaviour described above, we get

$$\begin{aligned} \mathbb{E}[Z^2 \mathbb{1}_{\{Z^2 > t^2\}}] &= t^2 \Pr(Z^2 > t^2) + \int_{t^2}^\infty \Pr(Z^2 > u) du \\ &\leq 2t^2 \left(e^{-\frac{c_3 t^2}{\sigma'^2}} + e^{-\frac{c_3 t \sqrt{d}}{\sigma'}} + e^{-\frac{c_3 t^2 d}{\sigma^2}} \right) + 2 \int_{t^2}^\infty e^{-\frac{c_3 u}{\sigma'^2}} du + 2 \int_{t^2}^\infty e^{-\frac{c_3 u d}{\sigma^2}} du \\ &\quad + 2 \int_{t^2}^\infty e^{-\frac{c_3 \sqrt{u d}}{\sigma'}} du \\ &\leq 2 \left(\frac{\sigma'^2}{c_3} + t^2 \right) e^{-\frac{c_3 t^2}{\sigma'^2}} + 2 \left(\frac{\sigma^2}{dc_3} + t^2 \right) e^{-\frac{c_3 t^2 d}{\sigma^2}} + 2 \left(\frac{2\sigma'^2}{dc_3^2} + \frac{2\sigma' t}{c_3 \sqrt{d}} + t^2 \right) e^{-\frac{c_3 t \sqrt{d}}{\sigma'}}. \end{aligned}$$

\square

Bounds on $\alpha'(\mathcal{M})$ and $\beta'(\mathcal{M})$ Recall that $Q_{\text{RM},j}$ denotes the rotated modulo quantizer without any subsampling for client $j \in \mathcal{C}_2$. From the description of RMQ in Algorithm 5, we have

$$Q_{\text{M},R_j,j}(\hat{g}_j(x), Y) = R_j^{-1} \left(\sum_{i \in [d]} Q_{\text{M}}(R_j \hat{g}_j(x)(i), R_j Y(i)) \cdot e_i \right).$$

The key step of the proof is bounding MSE and bias of RMQ. Towards that, we have the following lemma.

Remark 12. The calculation for MSE and bias are different in the proof of Lemma 9.6 compared to those in Proof of Lemma 3.2. This is because of the weaker concentration results available in this case.

Lemma 9.11. *Under the condition that $nr \geq c_4 d^2 \log(B/\sigma)$, we have for all $x \in \mathbb{R}^d$, $j \in \mathcal{C}_2$, and for some parameter $\delta \in (0, \sigma/\sqrt{c_3})$ that*

$$\begin{aligned} \mathbb{E} [\|Q_{\text{M},R_j,j}(\hat{g}_j(x), Y) - \hat{g}_j(x)\|_2^2] &\leq \frac{36\sigma^2}{c_3(\ell_2 - 2)^2} \left(\ln \frac{\sigma}{\sqrt{c_3}\delta} \right)^2 + 237\delta^2, \\ \mathbb{E} [Q_{\text{M},R_j,j}(\hat{g}_j(x), Y)] - \hat{g}_j(x) &\|_2^2 \leq 237\delta^2, \end{aligned}$$

where c_3 and c_4 are the universal constants same as in Lemma 9.6.

Proof. By considering events $\{|R_j(\hat{g}_j(x) - Y)(i)| \leq \Delta'\}$ and $\{|R_j(\hat{g}_j(x) - Y)(i)| \geq \Delta'\}$, and then using the facts for modulo quantizer, we have

$$\begin{aligned} \mathbb{E} [\|Q_{\text{M},R_j,j}(\hat{g}_j(x), Y) - \hat{g}_j(x)\|_2^2] &\leq d\varepsilon^2 + \sum_{i=1}^d \mathbb{E} \left[(Q_{\text{M},R_j,j}(\hat{g}_j(x), Y) - \hat{g}_j(x))(i)^2 \mathbb{1}_{\{|R_j(\hat{g}_j(x) - Y)(i)| \geq \Delta'\}} \right] \\ &\leq d\varepsilon^2 + 2\ell_2^2 \varepsilon^2 \sum_{i=1}^d \Pr(|R_j(\hat{g}_j(x) - Y)(i)| \geq \Delta') \\ &\quad + 2 \sum_{i=1}^d \mathbb{E} [(R_j(\hat{g}_j(x) - Y)(i))^2 \mathbb{1}_{\{|R_j(\hat{g}_j(x) - Y)(i)| \geq \Delta'\}}] \\ &\leq d\varepsilon^2 + 4\ell_2^2 \varepsilon^2 d \left(e^{-\frac{c_3 \Delta'^2}{\sigma'^2}} + e^{-\frac{c_3 \Delta' \sqrt{d}}{\sigma'}} + e^{-\frac{c_3 \Delta'^2 d}{\sigma'^2}} \right) \\ &\quad + 4d \left(\frac{\sigma'^2}{c_3} + \Delta'^2 \right) e^{-\frac{c_3 \Delta'^2}{\sigma'^2}} \\ &\quad + 4d \left(\frac{\sigma^2}{dc_3} + \Delta'^2 \right) e^{-\frac{c_3 \Delta'^2 d}{\sigma'^2}} + 4d \left(\frac{2\sigma'^2}{dc_3^2} + \frac{2\sigma' \Delta'}{c_3 \sqrt{d}} + \Delta'^2 \right) e^{-\frac{c_3 \Delta' \sqrt{d}}{\sigma'}}. \end{aligned} \tag{47}$$

For some parameter $\delta \in (0, \sqrt{d}\Delta)$, we substitute the other parameters as $\varepsilon = 2\Delta'/(\ell_2 - 2)$ and $\Delta'^2 = 9\Delta^2(\ln(\sqrt{d}\Delta/\delta))^2$, where $\Delta^2 = \max\{\frac{\sigma'^2}{c_3}, \frac{\sigma'^2}{dc_3^2}, \frac{\sigma^2}{dc_3}\}$. That gives

$$4\ell_2^2 \varepsilon^2 d \left(e^{-\frac{c_3 \Delta'^2}{\sigma'^2}} + e^{-\frac{c_3 \Delta' \sqrt{d}}{\sigma'}} + e^{-\frac{c_3 \Delta'^2 d}{\sigma'^2}} \right) = \left(\frac{12\ell_2}{\ell_2 - 2} \right)^2 d\Delta^2 \left(\ln \frac{\sqrt{d}\Delta}{\delta} \right)^2 \left(\frac{2}{(\sqrt{d}\Delta/\delta)^9} + \frac{1}{(\sqrt{d}\Delta/\delta)^3} \right)$$

$$\begin{aligned}
&\leq \left(\frac{12\ell_2}{\ell_2-2}\right)^2 d\Delta^2 \left(\ln \frac{\sqrt{d}\Delta}{\delta}\right)^2 \left(\frac{3}{(\sqrt{d}\Delta/\delta)^3}\right) \\
&= \left(\frac{12\ell_2}{\ell_2-2}\right)^2 \frac{\left(\ln \frac{\sqrt{d}\Delta}{\delta}\right)^2}{\sqrt{d}\Delta/\delta} \delta^2 \\
&\leq \left(\frac{24\ell_2}{e(\ell_2-2)}\right)^2 \delta^2 \\
&\leq 139\delta^2, \tag{48}
\end{aligned}$$

where the first inequality uses the fact that $\delta \in (0, \sqrt{d}\Delta)$, the second inequality uses $\ln x \leq 2\sqrt{x}/e$, and the final inequality uses the assumption that $n \geq 8$. For the last three terms in (47), we have

$$\begin{aligned}
&4d \left(\frac{\sigma'^2}{c_3} + \Delta'^2\right) e^{-\frac{c_3\Delta'^2}{\sigma'^2}} + 4d \left(\frac{\sigma^2}{dc_3} + \Delta'^2\right) e^{-\frac{c_3\Delta'^2 d}{\sigma^2}} + 4d \left(\frac{2\sigma'^2}{dc_3^2} + \frac{2\sigma'\Delta'}{c_3\sqrt{d}} + \Delta'^2\right) e^{-\frac{c_3\Delta'\sqrt{d}}{\sigma'}} \\
&\leq 8d(\Delta^2 + \Delta'^2) e^{-\frac{\Delta'^2}{\Delta^2}} + 4d(2\Delta^2 + 2\Delta\Delta' + \Delta'^2) e^{-\frac{\Delta'}{\Delta}} \\
&\leq 8d(\Delta^2 + \Delta'^2) e^{-\frac{\Delta'^2}{\Delta^2}} + 4d(3\Delta^2 + 2\Delta'^2) e^{-\frac{\Delta'}{\Delta}} \\
&= \frac{8d(\Delta^2 + 9\Delta^2(\ln(\sqrt{d}\Delta/\delta))^2)}{(\sqrt{d}\Delta/\delta)^9} + \frac{4d(3\Delta^2 + 18\Delta^2(\ln(\sqrt{d}\Delta/\delta))^2)}{(\sqrt{d}\Delta/\delta)^3} \\
&\leq \frac{4d(5\Delta^2 + 36\Delta^2(\ln(\sqrt{d}\Delta/\delta))^2)}{(\sqrt{d}\Delta/\delta)^3} \\
&\leq \frac{4d(5\Delta^2 + 144\Delta^2\sqrt{d}\Delta/(e^2\delta))}{(\sqrt{d}\Delta/\delta)^3} \\
&= \frac{20\delta^2}{(\sqrt{d}\Delta/\delta)} + \frac{576\delta^2}{e^2} \\
&\leq 98\delta^2, \tag{49}
\end{aligned}$$

where the first inequality is due to choice of Δ , the second is AM-GM inequality, the third one uses the fact: $\delta \in (0, \sqrt{d}\Delta)$, and the fourth one uses $\ln x \leq 2\sqrt{x}/e$. Substituting (48) and (49) in (47), we get

$$\mathbb{E} [\|Q_{\mathbf{M}, R_j, j}(\hat{g}_j(x), Y) - \hat{g}_j(x)\|_2^2] \leq \frac{36d\Delta^2}{(\ell_2-2)^2} \left(\ln \frac{\sqrt{d}\Delta}{\delta}\right)^2 + 237\delta^2.$$

Finally, we note that whenever $nr \geq c_4d^2 \log(B/\sigma)$, $\sigma'^2 \leq \sigma^2/d$ (see (46)). With our earlier choice of $\Delta^2 = \max\{\frac{\sigma'^2}{c_3}, \frac{\sigma'^2}{dc_3^2}, \frac{\sigma^2}{dc_3}\}$, this further implies $\Delta \leq \frac{\sigma}{\sqrt{dc_3}}$. Using this fact in the bound above establishes the MSE bound.

Bound for Bias. Using the fact that for $\{|R_j(\hat{g}_j(x) - Y)(i)| \leq \Delta'\}$ the modulo quantizer gives an unbiased estimate and the Jensen's inequality and, we have

$$\begin{aligned} \|\mathbb{E}[Q_{M,R_j,j}(\hat{g}_j(x), Y)] - \hat{g}_j(x)\|_2^2 &= \sum_{i=1}^d \mathbb{E}[(\hat{x}_R(i) - Rx(i)) \mathbb{1}_{|R(x-y)_i| \geq \Delta'}]^2 \\ &\leq \sum_{i=1}^d \mathbb{E}[(\hat{x}_R(i) - Rx(i))^2 \mathbb{1}_{|R(x-y)_i| \geq \Delta'}] \\ &\leq 237\delta^2. \end{aligned}$$

□

Completing the proof. We now calculate the MSE and bias for our Wyner-Ziv quantizer with subsampled RMQ.

Note that the inequality (37) derived in Lemma 3.3 holds in this case. Therefore, we have for $j \in \mathcal{C}_2$ that

$$\begin{aligned} \mathbb{E}[\|Q_{WZ,j}(\hat{g}_j(x), Y) - \hat{g}_j(x)\|_2^2] &\leq \frac{2d}{r_2} \mathbb{E}[\|Q_{M,R_j,j}(\hat{g}_j(x), Y) - \hat{g}_j(x)\|_2^2] \\ &\quad + \frac{2d}{r_2} \mathbb{E}[\|R\hat{g}_j(x) - RY\|_2^2] \\ &\leq \frac{2d}{r_2} \left(\frac{36\sigma^2}{c_3(\ell_2 - 2)^2} \left(\ln \frac{\sigma}{\sqrt{c_3}\delta} \right)^2 + 237\delta^2 \right) \\ &\quad + \frac{2d}{r_2} \mathbb{E}[\|R\hat{g}_j(x) - R\nabla f(x)\|_2^2] + \frac{2d}{r_2} \mathbb{E}[\|RY - R\nabla f(x)\|_2^2] \\ &\leq \frac{2d}{r_2} \left(\frac{36\sigma^2}{c_3(\ell_2 - 2)^2} \left(\ln \frac{\sigma}{\sqrt{c_3}\delta} \right)^2 + 237\delta^2 + \sigma^2 \right) \\ &\quad + \frac{2d}{r_2} \mathbb{E}[\|Y - \nabla f(x)\|_2^2], \end{aligned}$$

where the second last inequality uses bound from Lemma 9.6 and the fact

$$\mathbb{E}[\|R\hat{g}_j(x) - RY\|_2^2] = \mathbb{E}[\|R\hat{g}_j(x) - R\nabla f(x)\|_2^2] + \mathbb{E}[\|RY - R\nabla f(x)\|_2^2],$$

and the final inequality uses the fact that R is a unitary matrix and assumption (20).

Recall from Lemma 9.7 that the quantity $(\nabla f(x)(i) - Y(i))$ is subgaussian with variance parameter σ'^2 . Thus, $\mathbb{E}[(Y(i) - \nabla f(x)(i))^2] \leq c_5\sigma'^2$ for some universal constant $c_5 > 0$ (see, for instance, [56]). Again using the fact that $d\sigma'^2 < \sigma^2$, whenever $nr \geq c_4d^2 \log(B/\sigma)$, $\mathbb{E}[\|Y - \nabla f(x)\|_2^2] \leq \sigma^2$.

Further, the bias remains unchanged compared to without subsampling case, i.e.,

$$\|\mathbb{E}[Q_{WZ,j}(\hat{g}_j(x), Y)] - \hat{g}_j(x)\| = \|\mathbb{E}[Q_{M,R_j,j}(\hat{g}_j(x), Y)] - \hat{g}_j(x)\| \leq \sqrt{237}\delta.$$

At last, we set the following parameters for WZ-SGD:

$$\log \ell_2 = \lceil c \log \log nT \rceil, \quad \delta = \frac{2\sigma}{nT}.$$

Accordingly, we need to sample $r_2 = r / \lceil c \log \log nT \rceil$ coordinates at each client. Using the standard bounds for averaging of vectors in Lemma 2.1, respectively, we obtain

$$\alpha'^2(\mathcal{M}_t) \leq c_1 \frac{\sigma^2 \log \log nT}{n} \cdot \frac{d}{r}, \quad \beta'^2(\mathcal{M}_t) \leq \frac{c_2 \sigma^2}{n^2 T^2},$$

for suitably chosen constants $c_1, c_2 > 0$ and \mathcal{M}_t as defined in (26). The proof is completed by using the bounds on α' and β' with Lemma 5.2. \square

9.13 Proof of Theorem 5.5

We proceed in a way similar to the proof of Theorem 5.3. Towards that, we first use the mean square assumption in (27) to write

$$\mathbb{E} \left[\left\| \frac{2}{n} \sum_{i \in \mathcal{C}_2} \hat{g}_i(x_t) - \nabla f(x_t) \right\|^2 \right] \leq \frac{2\sigma^2}{n}.$$

Then, it only remains to bound the term $\mathbb{E} \left[\left\| \frac{2}{n} \sum_{i \in \mathcal{C}_2} Q_{\text{RDAQ}}(\hat{g}_i(x_t), Y_t) - \frac{2}{n} \sum_{i \in \mathcal{C}_2} \hat{g}_i(x_t) \right\|^2 \right]$, where $Y_t = \frac{2}{n} \sum_{i \in \mathcal{C}_2} Q_{\text{RATQ}}(\hat{g}_i(x_t))$ is the side-information. For that, we follow the proof of Lemma 4.3 in Section 9.8.

Fix any arbitrary client $i \in \mathcal{C}_2$. Conditioning on its gradient estimate $\hat{g}_i(x_t)$, the available side information at server $Y_t = y_t$, we have using the proof of Lemma 4.3 that

$$\mathbb{E} \left[\left\| Q_{\text{RDAQ}}(\hat{g}_i(x_t), Y_t) - \hat{g}_i(x_t) \right\|^2 \mid \hat{g}_i(x_t), Y_t = y_t \right] \leq \frac{16\sqrt{3}dB \|\hat{g}_i(x_t) - y_t\|}{\frac{r}{\lceil h + \log h \rceil} - 1}.$$

By the law of total expectations, we also have

$$\mathbb{E} \left[\left\| Q_{\text{RDAQ}}(\hat{g}_i(x_t), Y_t) - \hat{g}_i(x_t) \right\|^2 \right] \leq \frac{16\sqrt{3}dB \mathbb{E} [\|\hat{g}_i(x_t) - Y_t\|]}{\frac{r}{\lceil h + \log h \rceil} - 1}.$$

The proof is completed by noting that

$$\begin{aligned} \mathbb{E} [\|\hat{g}_i(x_t) - Y_t\|^2] &\leq \mathbb{E} [\|\hat{g}_i(x_t) - \nabla f(x_t)\|^2] + \mathbb{E} [\|\nabla f(x_t) - Y_t\|^2] \\ &\leq \sigma^2 + \frac{2\sigma^2}{n} + \frac{2dB^2}{n \left(\frac{r}{\lceil 3 + \log(1 + \ln^*(d/3)) \rceil} - 1 \right)}, \end{aligned}$$

where the first line is using Jensen's inequality, the only identity is due to the unbiased property of subsampled RATQ (*c.f.* Lemma 9.3), and the last line is due to (27) and applying the value of $\alpha(\mathcal{M}_t)$ for \mathcal{C}_1 in Theorem 5.3. \square

9.14 Proof of Theorem 6.1

The proof of this Theorem is similar to that of Lemma 3.2. We denote by $Q(X(i), Y(i))$ the output of the modulo quantizer with side information $Y(i)$ and parameters k, Δ' set as in (28). Then, we have

$$\begin{aligned} \mathbb{E} [\|Q_d(X, Y) - X\|^2] &\leq \sum_{i=1}^d \mathbb{E} [(Q(X(i), Y(i)) - X(i))^2] \\ &\leq \sum_{i=1}^d \mathbb{E} [(Q(X(i), Y(i)) - X(i))^2 \mathbb{1}_{\{|X(i) - Y(i)| \leq \Delta'\}}] \\ &\quad + \sum_{i=1}^d \mathbb{E} [(Q(X(i), Y(i)) - X(i))^2 \mathbb{1}_{\{|X(i) - Y(i)| \geq \Delta'\}}]. \end{aligned} \quad (50)$$

We bound the first term on the right-side in a similar manner as the bound in (34). Specifically, under the event $\{|X(i) - Y(i)| \leq \Delta'\}$, we get by Lemma 3.1 that

$$|Y(i) - X(i)| \leq \varepsilon = \frac{2\Delta'}{k-2}, \quad \text{almost surely,}$$

whereby

$$\sum_{i=1}^d \mathbb{E} [(Y(i) - X(i))^2 \mathbb{1}_{\{|X(i) - Y(i)| \leq \Delta'\}}] \leq d\varepsilon^2. \quad (51)$$

For the second term in the RHS note that $X(i) - Y(i)$ is subgaussian with variance factor σ_z^2 . Therefore, by proceeding in a similar manner as the derivation of (35) we get

$$\begin{aligned} &\sum_{i=1}^d \mathbb{E} [(Q(X(i), Y(i)) - X(i))^2 \mathbb{1}_{\{|X(i) - Y(i)| \geq \Delta'\}}] \\ &\leq 2 \sum_{i=1}^d [\mathbb{E} [(Q(X(i), Y(i)) - Y(i))^2 \mathbb{1}_{\{|X(i) - Y(i)| \geq \Delta'\}}] + \mathbb{E} [(Y(i) - X(i))^2 \mathbb{1}_{\{|X(i) - Y(i)| \geq \Delta'\}}]] \\ &\leq 2k^2\varepsilon^2 \sum_{i=1}^d P(|X(i) - Y(i)| \geq \Delta') + 2 \sum_{i=1}^d \mathbb{E} [(X(i) - Y(i))^2 \mathbb{1}_{\{|X(i) - Y(i)| \geq \Delta'\}}] \\ &\leq 4dk^2\varepsilon^2 e^{-d\Delta'^2/2\sigma_z^2} + 2 \sum_{i=1}^d \mathbb{E} [(X(i) - Y(i))^2 \mathbb{1}_{\{|X(i) - Y(i)| \geq \Delta'\}}] \\ &\leq 4dk^2\varepsilon^2 e^{-\Delta'^2/2\sigma_z^2} + 4(2\sigma_z^2 + d\Delta'^2) e^{-\frac{\Delta'^2}{2\sigma_z^2}}, \end{aligned} \quad (52)$$

where the second inequality follows upon noting from the description decoder of MQ in Alg. 3 that $|Q(X(i), Y(i)) - Y(i)| \leq \varepsilon k$ almost surely for each $i \in [d]$; the third inequality uses the fact that $X(i) - Y(i)$ is sub-Gaussian with variance parameter σ_z^2 ; and the fourth inequality is by Lemma 9.1.

Upon bounding the two terms on the right-side of (50) from above using (51), (52), we obtain

$$\mathbb{E} [\|Q_d(X, Y) - X\|^2] \leq d\varepsilon^2 + 4dk^2\varepsilon^2 e^{-\Delta'^2/2\sigma_z^2} + 4(2\sigma_z^2 + d\Delta'^2)e^{-\frac{\Delta'^2}{2\sigma_z^2}}.$$

Note that the RHS in the upper bound above is precisely the same as in (36) with σ_z^2 replacing Δ^2/d . Therefore proceeding in the same manner as in (36), we get

$$\mathbb{E} [\|Q_d(X, Y) - X\|^2] \leq 24 \frac{\sigma_z^2}{(k-2)^2} \ln \frac{\sigma_z}{\delta} + 154\delta^2.$$

Substituting the value of k and δ completes the proof. \square

9.15 Proof of Lemma 7.2

For $Q(x)$ as in (30), we have

$$Q(x) = \sum_{i=1}^N q_i/N,$$

where q_i for all $i \in \{1, \dots, N\}$ is an unbiased estimate of x and equals in distribution the output of the RDAQ quantizer for an input x and side information y . Moreover, q_i s are mutually independent conditioned on R . Therefore,

$$\begin{aligned} \mathbb{E} [\|Q(x) - x\|_2^2] &= \mathbb{E} \left[\left\| \sum_{i=1}^N \frac{q_i}{N} - x \right\|_2^2 \right] \\ &= \mathbb{E} \left[\mathbb{E} \left[\left\| \sum_{i=1}^N \frac{q_i}{N} - x \right\|_2^2 \middle| R \right] \right] \\ &= \mathbb{E} \left[\sum_{i=1}^N \frac{1}{N^2} \mathbb{E} [\|q_i - x\|_2^2 | R] \right] \\ &\leq 16\sqrt{3} \frac{\Delta}{N}, \end{aligned}$$

where the third identity follows from the conditional independence of q_i s after conditioning on R and the fact that q_i is an unbiased estimate of x . The final inequality follows from the fact that q_i equals in distribution the output of the RDAQ quantizer and then using Lemma 4.2. \square

References

- [1] J. Acharya, C. L. Canonne, P. Mayekar, and H. Tyagi, "Information-constrained optimization: can adaptive processing of gradients help?" *Advances in Neural Information Processing Systems*, 2021.
- [2] J. Acharya, C. L. Canonne, Z. Sun, and H. Tyagi, "Unified lower bounds for interactive high-dimensional estimation under information constraints," <http://arxiv.org/abs/2010.06562v5>, 2020.

- [3] J. Acharya, C. L. Canonne, P. Mayekar, and H. Tyagi, “Information-constrained optimization: can adaptive processing of gradients help?” <https://arxiv.org/abs/2104.00979>, 2021.
- [4] J. Acharya, C. De Sa, D. J. Foster, and K. Sridharan, “Distributed Learning with Sublinear Communication,” *International Conference on Machine Learning*, 2019.
- [5] A. Agarwal, P. L. Bartlett, P. Ravikumar, and M. J. Wainwright, “Information-Theoretic Lower Bounds on the Oracle Complexity of Stochastic Convex Optimization,” *IEEE Transactions on Information Theory*, vol. 5, no. 58, pp. 3235–3249, 2012.
- [6] N. Ailon and B. Chazelle, “Approximate nearest neighbors and the fast johnson-lindenstrauss transform,” *Proceedings of the ACM symposium on Theory of computing (STOC’06)*, pp. 557–563, 2006.
- [7] A. Albasyoni, M. Safaryan, L. Condat, and P. Richtárik, “Optimal gradient compression for distributed and federated learning,” *arXiv preprint arXiv:2010.03246*, 2020.
- [8] D. Alistarh, D. Grubic, J. Li, R. Tomioka, and M. Vojnovic, “QSGD: Communication-efficient SGD via gradient quantization and encoding,” *Advances in Neural Information Processing Systems*, pp. 1709–1720, 2017.
- [9] S. Augenstein, A. Hard, L. Ning, K. Singhal, S. Kale, K. Partridge, and R. Mathews, “Mixed federated learning: Joint decentralized and centralized learning,” *arXiv preprint arXiv:2205.13655*, 2022.
- [10] D. Basu, D. Data, C. Karakus, and S. Diggavi, “Qsparse-local-SGD: Distributed SGD with Quantization, Sparsification, and Local Computations,” *Advances in Neural Information Processing Systems*, 2019.
- [11] L. Bottou, “Large-scale machine learning with stochastic gradient descent,” *Proceedings of COMPSTAT’2010*, 2010.
- [12] S. Boucheron, G. Lugosi, and P. Massart, *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- [13] S. Bubeck, “Convex optimization: Algorithms and complexity,” *Foundations and Trends® in Machine Learning*, vol. 8, no. 3-4, pp. 231–357, 2015.
- [14] W.-N. Chen, P. Kairouz, and A. Özgür, “Breaking the communication-privacy-accuracy trilemma,” *Neural Information Processing Systems (NeurIPS)*, 2020.
- [15] P. Davies, V. Gurunathan, N. Moshrefi, S. Ashkboos, and D. Alistarh, “Distributed variance reduction with optimal communication,” *arXiv e-prints*, pp. arXiv–2002, 2020.
- [16] J. C. Duchi, M. I. Jordan, M. J. Wainwright, and Y. Zhang, “Optimality guarantees for distributed statistical estimation,” *arXiv:1405.0782*, 2014.
- [17] F. Faghri, I. Tabrizian, I. Markov, D. Alistarh, D. Roy, and A. Ramezani-Kebrya, “Adaptive gradient quantization for data-parallel sgd,” *Advances in Neural Information Processing Systems*, 2020.

- [18] G. D. Forney, “Coset codes. i. introduction and geometrical classification,” *IEEE Transactions on Information Theory*, vol. 34, no. 5, pp. 1123–1151, 1988.
- [19] V. Gandikota, D. Kane, R. Kumar Maity, and A. Mazumdar, “vqsgd: Vector quantized stochastic gradient descent,” in *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, ser. Proceedings of Machine Learning Research. PMLR, 2021, pp. 2197–2205.
- [20] A. Ghosh, R. K. Maity, and A. Mazumdar, “Distributed newton can communicate less and resist byzantine workers,” *Advances in Neural Information Processing Systems*, 2020.
- [21] T. Holenstein, “Parallel repetition: Simplification and the no-signaling case,” *Theory of Computing*, vol. 5, no. 8, pp. 141–172, 2009.
- [22] K. J. Horadam, *Hadamard matrices and their applications*. Princeton university press, 2012.
- [23] Z. Huang, W. Yilei, K. Yi *et al.*, “Optimal sparsity-sensitive bounds for distributed mean estimation,” *Advances in Neural Information Processing Systems*, pp. 6371–6381, 2019.
- [24] S. K. Jha, P. Mayekar, and H. Tyagi, “Fundamental limits of over-the-air optimization: Are analog schemes optimal?” *IEEE Journal on Selected Areas in Information Theory*, 2022.
- [25] D. Jhunjhunwala, A. Gadhikar, G. Joshi, and Y. C. Eldar, “Adaptive quantization of model updates for communication-efficient federated learning,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 3110–3114.
- [26] D. Jhunjhunwala, A. Mallick, A. H. Gadhikar, S. Kadhe, and G. Joshi, “Leveraging spatial and temporal correlations in sparsified mean estimation,” in *Advances in Neural Information Processing Systems*, 2021.
- [27] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings *et al.*, “Advances and open problems in federated learning,” *arXiv preprint arXiv:1912.04977*, 2019.
- [28] B. Kashin, “Section of some finite-dimensional sets and classes of smooth functions (in russian) izv,” *Acad. Nauk. SSSR*, vol. 41, pp. 334–351, 1977.
- [29] J. Konečný and P. Richtárik, “Randomized distributed mean estimation: Accuracy vs. communication,” *Frontiers in Applied Mathematics and Statistics*, vol. 4, p. 62, 2018.
- [30] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtarik, A. T. Suresh, and D. Bacon, “Federated learning: Strategies for improving communication efficiency,” *NIPS Workshop on Private Multi-Party Machine Learning*, 2016.
- [31] S. B. Korada and R. L. Urbanke, “Polar codes are optimal for lossy source coding,” *IEEE Transactions on Information Theory*, vol. 56, no. 4, pp. 1751–1768, 2010.
- [32] K. Liang and Y. Wu, “Improved communication efficiency for distributed mean estimation with side information,” in *2021 IEEE International Symposium on Information Theory (ISIT)*, 2021, pp. 3185–3190.

- [33] K. Liang, H. Zhong, H. Chen, and Y. Wu, “Wyner-Ziv Gradient Compression for Federated Learning,” <https://arxiv.org/abs/2111.08277>, 2021.
- [34] C.-Y. Lin, V. Kostina, and B. Hassibi, “Differentially Quantized Gradient Descent,” in *IEEE International Symposium on Information Theory (ISIT)*, 2021.
- [35] C. Ling, S. Gao, and J. Belfiore, “Wyner-ziv coding based on multidimensional nested lattices,” *IEEE Transactions on Communications*, vol. 60, no. 5, pp. 1328–1335, 2012.
- [36] L. Liu, “Polar codes and polar lattices for efficient communication and source quantization,” *Ph.D. Thesis*, 2016.
- [37] L. Liu and C. Ling, “Polar lattices are good for lossy compression,” *CoRR*, vol. abs/1501.05683, 2015.
- [38] Y. Lu and C. De Sa, “Moniqua: Modulo quantized communication in decentralized sgd,” *arXiv preprint arXiv:2002.11787*, 2020.
- [39] Y. Lyubarskii and R. Vershynin, “Uncertainty principles and vector quantization,” *IEEE Transactions on Information Theory*, vol. 56, no. 7, pp. 3491–3501, 2010.
- [40] P. Mayekar, S. K. Jha, and H. Tyagi, “Wyner-ziv compression is (almost) optimal for distributed optimization,” in *2022 IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2022, pp. 578–583.
- [41] P. Mayekar, A. T. Suresh, and H. Tyagi, “Wyner-Ziv estimators: Efficient distributed mean estimation with side-information,” in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2021, pp. 3502–3510.
- [42] P. Mayekar and H. Tyagi, “Limits on gradient compression for stochastic optimization,” *Proceedings of the IEEE International Symposium of Information Theory (ISIT’ 20)*, 2020.
- [43] ———, “RATQ: A universal fixed-length quantizer for stochastic optimization,” *IEEE Transactions on Information Theory*, 2020.
- [44] A. Nemirovsky, “Information-based complexity of convex programming,” 1995, Available Online http://www2.isye.gatech.edu/ne-mirovs/Lec_EMCO.pdf.
- [45] A. Nemirovsky and D. B. Yudin, “Problem complexity and method efficiency in optimization.” *Wiley series in Discrete Mathematics and Optimization*, 1983.
- [46] Y. Oohama, “Gaussian multiterminal source coding,” *IEEE Transactions on Information Theory*, vol. 43, no. 6, pp. 1912–1923, 1997.
- [47] S. S. Pradhan and K. Ramchandran, “Distributed source coding using syndromes (discus): design and construction,” *IEEE Transactions on Information Theory*, vol. 49, no. 3, pp. 626–643, 2003.
- [48] A. Ramezani-Kebrya, F. Faghri, and D. M. Roy, “Nuqsgd: Improved communication efficiency for data-parallel sgd via nonuniform quantization,” *arXiv preprint arXiv:1908.06077*, 2019.

- [49] A. Reisizadeh, A. Mokhtari, H. Hassani, A. Jadbabaie, and R. Pedarsani, “Fedpaq: A communication-efficient federated learning method with periodic averaging and quantization,” in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2020, pp. 2021–2031.
- [50] M. Safaryan, E. Shulgin, and P. Richtárik, “Uncertainty principle for communication compression in distributed and federated learning and the search for an optimal compressor,” *arXiv preprint arXiv:2002.08958*, 2020.
- [51] R. Saha, S. Rini, M. Rao, and A. Goldsmith, “Decentralized optimization over noisy, rate-constrained networks: How we agree by talking about how we disagree,” in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 5055–5059.
- [52] F. Seide, H. Fu, J. Droppo, G. Li, and D. Yu, “1-bit stochastic gradient descent and its application to data-parallel distributed training of speech dnns,” *Fifteenth Annual Conference of the International Speech Communication Association*, 2014.
- [53] S. U. Stich, J.-B. Cordonnier, and M. Jaggi, “Sparsified sgd with memory,” *Advances in Neural Information Processing Systems 31*, 2018.
- [54] A. T. Suresh, Z. Sun, J. Ro, and F. Yu, “Correlated quantization for distributed mean estimation and optimization,” in *International Conference on Machine Learning*. PMLR, 2022, pp. 20 856–20 876.
- [55] A. T. Suresh, F. X. Yu, S. Kumar, and H. B. McMahan, “Distributed mean estimation with limited communication,” *Proceedings of the International Conference on Machine Learning (ICML’ 17)*, vol. 70, pp. 3329–3337, 2017.
- [56] R. Vershynin, *High-dimensional probability: An introduction with applications in data science*. Cambridge university press, 2018, vol. 47.
- [57] T. Vogels, S. P. Karimireddy, and M. Jaggi, “Powersgd: Practical low-rank gradient compression for distributed optimization,” 2019.
- [58] H. Wang, S. Sievert, S. Liu, Z. Charles, D. Papailiopoulos, and S. Wright, “Atomo: Communication-efficient learning via atomic sparsification,” *Advances in Neural Information Processing Systems*, pp. 9850–9861, 2018.
- [59] J. Wangni, J. Wang, J. Liu, and T. Zhang, “Gradient sparsification for communication-efficient distributed optimization,” *Advances in Neural Information Processing Systems*, 2018.
- [60] W. Wen, C. Xu, F. Yan, C. Wu, Y. Wang, Y. Chen, and H. Li, “TernGrad: Ternary gradients to reduce communication in distributed deep learning,” *Advances in Neural Information Processing Systems*, pp. 1509–1519, 2017.
- [61] A. Wyner and J. Ziv, “The rate-distortion function for source coding with side information at the decoder,” *IEEE Transactions on information Theory*, vol. 22, no. 1, pp. 1–10, 1976.
- [62] R. Zamir, S. Shamai, and U. Erez, “Nested linear/lattice codes for structured multiterminal binning,” *IEEE Transactions on Information Theory*, vol. 48, no. 6, pp. 1250–1276, 2002.