# Differentially Quantized Gradient Methods

Chung-Yi Lin, Victoria Kostina, and Babak Hassibi

Abstract—Consider the following distributed optimization scenario. A worker has access to training data that it uses to compute the gradients while a server decides when to stop iterative computation based on its target accuracy or delay constraints. The server receives all its information about the problem instance from the worker via a rate-limited noiseless communication channel.

We introduce the principle we call differential quantization (DQ) that prescribes compensating the past quantization errors to direct the descent trajectory of a quantized algorithm towards that of its unquantized counterpart. Assuming that the objective function is smooth and strongly convex, we prove that differentially quantized gradient descent (DQ-GD) attains a linear contraction factor of  $\max\{\sigma_{\rm GD}, \rho_n 2^{-R}\}$ , where  $\sigma_{\rm GD}$  is the contraction factor of unquantized gradient descent (GD),  $\rho_n \geq 1$  is the covering efficiency of the quantizer, and R is the bitrate per problem dimension n. Thus at any  $R \ge \log_2 \rho_n / \sigma_{\rm GD}$ bits, the contraction factor of DQ-GD is the same as that of unquantized GD, i.e., there is no loss due to quantization. We show a converse demonstrating that no algorithm within a certain class can converge faster than  $\max\{\sigma_{\mathrm{GD}}, 2^{-R}\}$ . Since quantizers exist with  $\rho_n \to 1$  as  $n \to \infty$  (Rogers, 1963), this means that DQ-GD is asymptotically optimal. In contrast, naively quantized GD where the worker directly quantizes the gradient barely attains  $\sigma_{\rm GD} + \rho_n 2^{-R}$ .

The principle of differential quantization continues to apply to gradient methods with momentum such as Nesterov's accelerated gradient descent, and Polyak's heavy ball method. For these algorithms as well, if the rate is above a certain threshold, there is no loss in contraction factor obtained by the differentially quantized algorithm compared to its unquantized counterpart, and furthermore, the differentially quantized heavy ball method attains the optimal contraction achievable among all (even unquantized) gradient methods.

Experimental results on least-squares problems validate our theoretical analysis.

Index Terms—gradient descent, quantized gradient descent, accelerated gradient descent, heavy ball method, error compensation, error feedback, sigma-delta modulation, federated learning, linear convergence.

# I. INTRODUCTION

#### A. Motivation and related work

Distributed optimization plays a central role in largescale machine learning where gradient descent (GD) and its stochastic variant SGD are employed to minimize an objective function [2]–[9]. Despite the scalability of parallel gradient training, the frequent exchange of high-dimensional gradients between distributed agents in the federated learning setting

Chung-Yi Lin (hsnu1220@gmail.com) is with Kronos Research, Taiwan. V. Kostina (vkostina@caltech.edu) and B. Hassibi (hassibi@caltech.edu) are with California Institute of Technology. This work was supported in part by the National Science Foundation (NSF) under grants CCF-1751356, CCF-1956386, CNS-0932428, CCF-1018927, CCF-1423663 and CCF-1409204, by a grant from Qualcomm Inc., by NASA's Jet Propulsion Laboratory through the President and Director's Fund, and by King Abdullah University of Science and Technology. This paper was presented in part at ISIT 2021 [1].

has become a communication bottleneck that slows down the overall learning process [3], [6], [10]–[13].

A natural approach to alleviating that communication bottleneck is to quantize the gradients with a limited number of bits per problem dimension. Its power was first demonstrated by Seide et al. [10], where the gradient computed by stochastic gradient descent (SGD) [14] is quantized down to just one bit per dimension and the quantization error is carried forward across mini-batches, resulting in almost no loss in empirical convergence performance compared to the unquantized algorithm. Wen at al. [15] propose a ternary quantizer for SGD and prove that it converges almost surely under the assumption of bounded gradients. Bernstein et al. [16] propose a signbased quantizer for mini-batch SGD, give its convergence analysis on nonconvex problems, and extend it to accelerated gradient descent and to a multi-worker setting. Alistarh et al. [17] propose a quantized SGD algorithm that compresses the gradient using a stochastic scalar quantizer with an adjustable number of quantization levels, and provide convergence guarantees that depend on this variable compression rate on smooth convex and non-convex functions. The quantizer in [10], [15]-[17] is a uniform scalar quantizer, which simply rounds the binary representation of each coordinate to a fixed number of bits, while [18] considers a non-uniform scalar quantizer, and [19], [20] construct vector quantizers from the convex hull of specifically structured point sets.

A different approach to addressing the communication bottleneck in parallel SGD training is to sparsify the gradient vectors [12], [21]–[25]. For example, the top-k sparsifier (or *compressor*) preserves the k coordinates of the largest magnitude and sends them with full precision [12], [21], [22], [26]–[28]. A user-specified parameter (e.g. k for the top-k compressor) serves as a proxy for the communication rate in this line of work.

For an empirical risk minimization problem where the global objective function is the average of local objective functions, recent works [29]–[31] perform analog gradient compression and communication by taking the physical superposition nature of the underlying multiple-access channel into the account.

The assumption of unbiased compression error [32]–[43] is commonly imposed to enable convergence analyses of compressed SGD. Employing biased compressors in compressed SGD can lead to divergence: for example, both the 1-bit SGD without mini-batching [10], [44] and the top-1 compressed SGD [45] diverge on some problem instances. A set of sufficient conditions on the compression operators to ensure convergence of SGD is put forth in [46].

The same paper - [10] - that initiates the study of quantized SGD is also the first to introduce the idea of adding back previous quantization errors before quantizing the gradient

at the next step of iterative optimization, which fixes the divergence issue mentioned above. The idea, referred to as error compensation, or error feedback, in the federated learning literature, has been long known as  $\Sigma$ - $\Delta$  modulation [47] in the information theory literature. Stich et al. [22] apply the mechanism of error feedback in [10] to a more general setting of SGD and show that it converges with the same order as unquantized SGD on strongly convex and smooth functions, providing the first theoretical performance guarantee of that error feedback strategy. Karimireddy et al. [44] extend the analysis of [22] to the non-convex and weakly convex objective functions, while Zheng et al. [48] and Gorbunov et al. [40] prove its convergence in the multi-worker setting. Wu et al. [49] propose an error feedback mechanism different from [10] and prove its convergence on quadratic functions using the same quantizer as in [17]. Past quantization errors in the algorithm of [49] accumulate from one iteration to another and are weighted by time-decaying factors. The momentum correction used in [25], as well as the distributed SGD with skipped communication rounds in [50], are also variants of error compensation. Qian et al. [51] propose an error-compensated accelerated SGD, while Richrárik et al. [52] propose an errorcompensated SGD that achieves the same order of convergence as SGD with unbiased compressors [42]. The analyses in [22], [40], [44], [52] assume that the compressor is a contraction operator, while [49] also assumes its unbiasedness. Horýath and Richtárik [53] construct an unbiased compressor from a contractive compressor and employ the resulting unbiased compressor within SGD as an alternative to error feedback to overcome the divergence issue with biased compressors.

Although a number of works provide convergence analyses of their proposed methods, showing that convergence rates of quantized gradient methods depend on the bit rate R [17]–[19], [54], [55], there are few existing convergence lower bounds in terms of R that apply to any algorithm within a specified class. For quantized projected SGD, [19], [55] give lower bounds to a minimax expected estimation error (i.e. difference between the output function value and the optimal one), which is in the same order of convergence as that of the unquantized SGD over convex functions. However, the allowable quantizer input in [19], [55] is fixed to be the gradient of the current iterate, precluding the use of error compensation.

The parameter server framework that we consider in this paper is somewhat different from the *distributed estimation* or optimization setting [56], where there has also been great interest in communication-efficient algorithms to account for the distributed nature of these problems. In such applications, all parties in a connected network communicate back and forth in order to estimate the mean of a distribution [57] (or a population [58]) or to solve a convex optimization cooperatively with quantization effects [59], [60]. Information-theoretic lower bounds have also been established either in the minimax sense for distributed statistical estimation problem [61] or in terms of the communication complexity for the distributed convex learning problem [58], [62], [63].

#### B. Contributions

In this paper, we provide a lower bound on quantized non-stochastic gradient descent, and we show a single-worker algorithm that achieves the lower bound with equality, thereby establishing an information-theoretic fundamental limit of quantized gradient descent. In other words, we quantify exactly the minimum amount of information required to achieve a desired convergence speed (within a class of algorithms), and we show an algorithm that achieves it. Because the algorithm achieves the information-theoretic converse with equality, no other algorithm can surpass its performance. It is remarkable that only a finite bit rate is required to achieve the optimal convergence speed achievable with an infinite rate. Our analysis is sharper than existing analyses because we identify constants and not just the order of convergence. We focus on (nonstochastic) GD and not on SGD as most prior work. We do not assume that the quantizer is unbiased or is a contraction operator - our information-theoretic lower bound applies to any quantizer, and any quantizer can be inserted into our algorithm, although our achievability result suggests that picking a (scalar or vector) quantizer with good covering efficiency would perform best. Our mechanism for error compensation (that we call "differential quantization") differs from prior works in the gradient-compute point, which is crucial for achieving our sharp information-theoretic lower bound. Although our information-theoretic lower bound applies to the multi-worker setting as well, our best achievability bound comes short of it at a finite R. Thus, it remains an open problem whether the lower bound is achievable in the multi-worker setting. While our main results are presented in terms of a quantity that is asymptotic in the number of iterations T, the analyses that lead to these results are nonasymptotic. Incidentally, we discover two new results on the classical (unquantized) gradient methods: a slightly more general converse for the gradient descent, and a nonasymptotic global convergence bound on Polyak's heavy ball method.

We consider the single-worker scenario of the parameter server framework [11], [15]–[18], [64], [65] consisting of a worker that computes the gradients and a server that successively refines the model parameter (i.e. the iterate) and decides when to stop the distributed iterative algorithm based on its target accuracy or delay constraints. See Fig. 1.

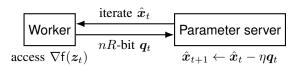


Fig. 1: Quantized gradient descent in a single-worker remote training setting. At each iteration t, the server first sends the current iterate  $\hat{x}_t$  to the worker noiselessly, who computes the gradient at some point  $z_t$  that is a function of (but not necessarily equal to)  $\hat{x}_t$ . Then, the worker forms a descent direction  $q_t$  and pushes it back to the server under the nR bits per iteration constraint.

We study the fundamental tradeoff between the convergence rate and the communication rate of quantized gradient descent. We focus on the class  $\mathcal{F}_n$  of smooth and strongly convex objective functions  $f \colon \mathbb{R}^n \mapsto \mathbb{R}$  whose minimizers are bounded in the Euclidean norm. For a quantized iterative algorithm A, its worst-case linear *contraction factor* over  $\mathcal{F}_n$  at rate R bits per problem dimension is defined as

$$\sigma_{\mathsf{A}}(n,R) \triangleq \inf_{R' \leq R} \sup_{\mathsf{f} \in \mathcal{F}_n} \limsup_{T \to \infty} \|\hat{\boldsymbol{x}}_T(R') - \boldsymbol{x}_\mathsf{f}^*\|^{\frac{1}{T}}$$
 (1)

where  $x_f^*$  is the optimizer, and  $\hat{x}_0(R'), \hat{x}_1(R'), \hat{x}_2(R'), \ldots$  is the sequence of iterates generated by A in response to  $f \in \mathcal{F}_n$  when it operates at R' bits per problem dimension.

We consider three popular algorithms that converge linearly: the classical gradient descent (GD) with fixed step size, the accelerated gradient descent (AGD) [66], and the heavy ball method (HB) [67]. We propose a principle for error feedback we call *differential quantization* (DQ) that says that the quantizer input should be formed in such a way as to guide the descent trajectory of the quantized algorithm towards the descent trajectory of its unquantized counterpart. By applying the DQ principle to the GD, AGD, and HB algorithms, we construct three new quantized iterative optimization algorithms: DQ-GD, DQ-AGD, and DQ-HB. By analyzing them, we show achievability bounds of the form<sup>2</sup>

$$\sigma_{\mathcal{A}}(n,R) \le \max \left\{ \sigma_{\mathcal{A}}, \rho_n 2^{-R} \phi_{\mathcal{A}}(n,R) \right\}, \tag{2}$$

where  $A \in \{DQ\text{-}GD, DQ\text{-}AGD, DQ\text{-}HB\}$ ,  $\sigma_A$  is the contraction factor of the unquantized counterpart of A,  $\rho_n \geq 1$  is the covering efficiency of the quantizer, and  $\phi_A(n,R) \geq 1$  is function that we specify; for example,

$$\sigma_{\text{DQ-GD}}(n, R) \le \max \left\{ \sigma_{\text{GD}}, \rho_n 2^{-R} \right\}.$$
 (3)

As (2) indicates, each of the novel DQ algorithms achieves the corresponding  $\sigma_A$  once the rate passes a hard threshold. In other words, there is no loss at all due to quantization once the rate is high enough.

We show an information-theoretic converse of the form

$$\sigma_{\mathcal{A}}(n,R) \ge \max\left\{\sigma_{\mathcal{GD}}, 2^{-R}\right\},\tag{4}$$

which applies to any "quantized gradient descent" algorithm A (in the class of "quantized gradient descent" algorithms, summarized in Fig. 1, the server can utilize only the last quantized input to form the next iterate). Recalling the classical result of Rogers [68, Th. 3] that shows the existence of quantizers with covering efficiency  $\rho_n \to 1$  as  $n \to \infty$  and comparing (2) and (4), one can deduce the asymptotic optimality of DQ-GD within the class of "GD-like" algorithms. In contrast, the natural method that quantizes the gradient of its current iterate directly [15]–[17], [69] referred to as naively quantized (NQ) GD in this paper, has contraction factor (in the single-worker scenario; see Section V for the multi-worker result)

$$\sigma_{\text{NQ-GD}}(n,R) \le \sigma_{\text{GD}} + \frac{2\kappa}{\kappa + 1} \rho_n 2^{-R}$$
 (5)

<sup>1</sup>The term "linear convergence" is used in the literature as a synonym for convergence with the rate of geometric progression. Note that SGD converges only sub-linearly over smooth and strongly convex functions [32]–[34].

<sup>2</sup>The convergence result on DQ-HB in (2) requires that the function  $f \in \mathcal{F}_n$  is twice continuously differentiable.

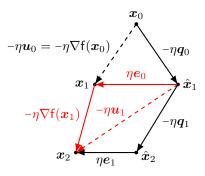


Fig. 2: Illustration of the DQ-GD algorithm (Algorithm 1).

where  $\kappa \geq 1$  is the condition number of f. The guarantee (5) is significantly worse than (3).

Our numerical results indicate that the upper bounds (2) and (5) accurately represent the actual achieved contraction factors.

Within a wider class of quantized gradient methods (the server can utilize full memory of the past), the converse (4) can be surpassed. Once the rate passes the threshold mentioned earlier, DQ-HB attains the minimum possible contraction factor among all algorithms in that wider class, even unquantized ones.

The rest of the paper is organized as follows. Differentially quantized algorithms are presented in Section II. Their convergence analyses and an experimental validation on least-squares problems are shown in Section III. The converses are presented in Section IV. The multi-worker setting is discussed in Section V.

## II. DIFFERENTIALLY QUANTIZED ALGORITHMS

#### A. Quantizers employed in DO algorithms

A quantizer of dimension n and rate R is a function  $q: \mathcal{D} \to \mathbb{R}^n$ , where  $\mathcal{D} \subseteq \mathbb{R}^n$  is the domain, such that the image of q satisfies

$$|\operatorname{Im}(\mathsf{q})| = 2^{nR}.\tag{6}$$

This is the classical general fixed-rate quantizer in the information theory literature. We fix a dimension-n, rate-R quantizer q, and we set up quantizer  $q_t$  to be used at iteration t as

$$q_t(\cdot) = r_t q(\cdot/r_t) \tag{7}$$

for a properly chosen sequence of shrinkage factors  $\{r_t\}$  (see (22), (35), and (46), below). Therefore, each quantizer  $q_t$  has the same geometric structure but different resolution.

## B. Differentially Quantized Gradient Descent

The (unquantized) gradient descent algorithm searches along the direction of the negative gradient toward which the function value decreases:

$$\boldsymbol{x}_{t+1} = \boldsymbol{x}_t - \eta \nabla f(\boldsymbol{x}_t), \tag{8}$$

where  $\eta > 0$  is the constant stepsize chosen to minimize the function value along the search direction.

In Fig. 2, we illustrate an application of differential quantization (DQ) to GD (8), which yields the DQ-GD algorithm

(Algorithm 1). At each iteration t, DQ-GD first determines the iterate  $x_t$  associated with the corresponding unquantized algorithm, i.e., GD, by compensating previous scaled quantization error  $\eta e_{t-1}$  (Line 4). It then computes the gradient at  $z_t = x_t$  (Lemma B.1) and sets the quantizer input as (Line 5)

$$\mathbf{u}_t = \nabla f(\hat{\mathbf{x}}_t + \eta \mathbf{e}_{t-1}) - \mathbf{e}_{t-1}, \tag{9}$$

which in the absence of quantization error  $e_t$  would guide the iterate  $\hat{x}_t$  back to  $x_{t+1}$  (see Fig. 2). The recorded scaled quantization error  $\eta e_t$  captures exactly the difference between  $\hat{x}_{t+1}$  and  $x_{t+1}$  for the next iteration.

See Appendix A for the DQ-GD algorithm with varying stepsize  $\eta_t$ .

# Algorithm 1: DQ-GD

```
1 Initialize e_{-1} = 0

2 for t = 0, 1, 2, \dots do

3 | Worker:

4 | z_t = \hat{x}_t + \eta e_{t-1}

5 | u_t = \nabla f(z_t) - e_{t-1}

6 | q_t = q_t(u_t)

7 | e_t = q_t - u_t

8 | Server: \hat{x}_{t+1} = \hat{x}_t - \eta q_t

9 end
```

#### C. Differentially Quantized Accelerated Gradient Descent

Nesterov's Accelerated Gradient Descent (AGD) [70] keeps track of two iterate sequences

$$\mathbf{y}_{t+1} = \mathbf{x}_t - \eta \nabla f(\mathbf{x}_t) \tag{10}$$

$$x_{t+1} = y_{t+1} + \gamma (y_{t+1} - y_t).$$
 (11)

It first performs the gradient descent step (10), and then adds the momentum term  $\gamma\left(y_{t+1}-y_{t}\right)$  (11) to form a projection  $x_{t+1}$  of the GD iterate  $y_{t+1}$  to its near future. The momentum term incorporates second-order effects by leveraging the past  $y_{t}$ . The AGD is the first algorithm that achieved the contraction factor that is order-wise optimal (in terms of the condition number of f) among all first-order (gradient) optimization methods [66] (Lemma B.5). There are various interpretations of Nesterov's acceleration phenomenon. We refer the reader to [71] for a connection to the mirror descent algorithm and to [72] for an interpretation in terms of differential equation.

Differentially Quantized AGD algorithm is presented as Algorithm 2. At each iteration t, DQ-AGD uses the past two quantization errors  $e_{t-1}, e_{t-2}$  to determine the gradient-compute point  $z_t$  (Line 4) and the quantizer input  $u_t$  (Line 5). As dictated by the principle of differential quantization, DQ-AGD computes the gradient at the same point as unquantized AGD, i.e.,  $z_t = x_t$  (Lemma B.4).

# D. Differentially Quantized Heavy Ball Method Polyak's Heavy Ball (HB) algorithm [67] iterates

$$\boldsymbol{x}_{t+1} = \boldsymbol{x}_t - \eta \nabla f(\boldsymbol{x}_t) + \gamma \left(\boldsymbol{x}_t - \boldsymbol{x}_{t-1}\right), \tag{12}$$

# Algorithm 2: DQ-AGD

```
1 Initialize e_{-2} = e_{-1} = \mathbf{0}, \ \hat{\pmb{y}}_0 = \hat{\pmb{x}}_0
 2 for t = 0, 1, 2, \dots do
               Worker:
                        z_t = \hat{x}_t + \eta [e_{t-1} + \gamma (e_{t-1} - e_{t-2})]
 4
                        \boldsymbol{u}_t = \nabla f(\boldsymbol{z}_t) - [\boldsymbol{e}_{t-1} + \gamma (\boldsymbol{e}_{t-1} - \boldsymbol{e}_{t-2})]
 5
                        q_t = q_t(u_t)
  6
 7
                        e_t = q_t - u_t
  8
                        \hat{\boldsymbol{y}}_{t+1} = \hat{\boldsymbol{x}}_t - \eta \boldsymbol{q}_t
 9
                        \hat{\boldsymbol{x}}_{t+1} = \hat{\boldsymbol{y}}_{t+1} + \gamma \left( \hat{\boldsymbol{y}}_{t+1} - \hat{\boldsymbol{y}}_{t} \right)
10
11 end
```

where  $\gamma(x_t - x_{t-1})$  is the momentum term that nudges  $x_{t+1}$  in the direction of the previous step, and accelerates convergence to the optimizer. In contrast to AGD, the HB method only uses the gradient at the current iterate. The HB method derives from the analogy with physics, since the continuous-time counterpart of (12) is a second-order ODE that models the motion of a body ("the heavy ball") in a field with potential f under the force of friction. At the expense of requiring function f in  $\mathcal{F}_n$  to be further twice continuously differentiable, the HB algorithm can be shown to converge with the optimal contraction factor achievable among all first-order optimization methods [67, Th. 3.1] (Lemma B.8), [66, Th. 2.1.13] (Lemma C.2). In comparison, the AGD approaches it only order-wise, but it does not require the second derivative of f to exist, a significant restriction in practical applications.

Differentially Quantized HB algorithm is presented as Algorithm 3. In accordance with the principle of differential quantization, the worker computes the gradient at  $z_t = x_t$  (Lemma B.7). Note that DQ-HB has the same expression for its quantizer input  $u_t$  (Line 4) as DQ-AGD (Line 5).

# **Algorithm 3:** Differentially Quantized Heavy Ball Method (DQ-HB)

```
1 Initialize e_{-2}=e_{-1}=0,\;\hat{x}_{-1}=\hat{x}_0
2 for t=0,1,2,\ldots do
3 | Worker:
4 | z_t=\hat{x}_t+\eta e_{t-1}
5 | u_t=\nabla \mathsf{f}(z_t)-[e_{t-1}+\gamma\,(e_{t-1}-e_{t-2})]
6 | q_t=\mathsf{q}_t(u_t)
7 | e_t=q_t-u_t
8 | Server: \hat{x}_{t+1}=\hat{x}_t-\eta q_t+\gamma\,(\hat{x}_t-\hat{x}_{t-1})
9 end
```

## III. CONVERGENCE RATES OF DQ ALGORITHMS

#### A. Definitions

We denote by  $\|\cdot\|$  the Euclidean norm, and by  $\mathcal{B}(r) \triangleq \{ \boldsymbol{u} \in \mathbb{R}^n \colon \|\boldsymbol{u}\| \leq r \}$  the Euclidean ball in  $\mathbb{R}^n$  with radius r and center at  $\boldsymbol{0}$ .

We fix positive scalars L, and  $\mu$ , and D, and we say that a continuously differentiable function  $f: \mathbb{R}^n \to \mathbb{R}$  is in class (12)  $\mathcal{F}_n$  if

i) f is L-smooth, i.e.,

$$\|\nabla f(\boldsymbol{v}) - \nabla f(\boldsymbol{w})\| \le L \|\boldsymbol{v} - \boldsymbol{w}\|; \tag{13}$$

ii) f is  $\mu$ -strongly convex, i.e.,

function 
$$v \mapsto f(v) - \frac{\mu}{2} \|v\|^2$$
 is convex; (14)

iii) the minimizer  $x_f^* \triangleq \arg\min_{x \in \mathbb{R}^n} f(x)$  satisfies

$$\|\boldsymbol{x}_{\mathsf{f}}^* - \hat{\boldsymbol{x}}_0\| \le D,\tag{15}$$

where  $\hat{x}_0$  is the starting location of iterative algorithms. We say that  $f \colon \mathbb{R}^n \mapsto \mathbb{R}$  is in class  $\mathcal{F}_n^2$  if it is in  $\mathcal{F}_n$  and is in addition twice continuously differentiable.

We denote the *condition number* of an  $f \in \mathcal{F}_n$  by

$$\kappa \triangleq \frac{L}{\mu}.\tag{16}$$

Note that  $\kappa \geq 1$  due to (13) and (14).

For a bounded-domain quantizer  $q: \mathcal{D} \to \mathbb{R}^n$ , we refer to

$$r(q) \triangleq \max \{ \delta \colon \mathcal{B}(\delta) \subseteq \mathcal{D} \}$$
 (17)

as the dynamic range of q, to

$$d(q) \triangleq \min \{d \colon \forall x \in \mathcal{D}, \ \|x - q(x)\| \le d\}$$
 (18)

as its covering radius, and to

$$\rho(\mathbf{q}) \triangleq \left| \operatorname{Im}(\mathbf{q}) \right|^{1/n} \frac{\mathsf{d}(\mathbf{q})}{r(\mathbf{q})} \tag{19}$$

as its covering efficiency.3 A scalar uniform quantizer qu has domain  $[-r(q_u), r(q_u)]^n$  and covering efficiency  $\sqrt{n}$ . This is wasteful: the classical result of Rogers [68, Th. 3] implies that there exists a sequence of n-dimensional quantizers  $q_n$  with  $\rho(q_n) \to 1$  as  $n \to \infty$ , while definition (19) implies that  $\rho(q) \ge 1$  for any quantizer q.

#### B. DO-GD: convergence analysis and simulation results

Unquantized gradient descent with the optimal stepsize given by

$$\eta = \eta_{\rm GD} \triangleq \frac{2}{L+\mu} \tag{20}$$

achieves contraction factor

$$\sigma_{\rm GD} \triangleq \frac{\kappa - 1}{\kappa + 1} \tag{21}$$

over  $\mathcal{F}_n$  [67, Th. 1.4], [66, Th. 2.1.15] (Lemma B.2). The following result provides a convergence guarantee for DQ-GD.

**Theorem III.1** (Convergence of DO-GD). Fix a dimension-n, rate-R quantizer q with dynamic range 1 and covering efficiency  $\rho_n$ . Then, Algorithm 1 with stepsize (20) and dynamic ranges  $r_0 = LD$ ,

$$r_{t+1} = \sigma_{\text{GD}}^{t+1} LD + r_t \rho_n 2^{-R}, \quad t = 1, 2, \dots$$
 (22)

in the definition of  $q_t$  (7) achieves the following contraction factor over  $\mathcal{F}_n$  (1):

$$\sigma_{\text{DQ-GD}}(n,R) \le \max\left\{\sigma_{\text{GD}}, \rho_n 2^{-R}\right\}.$$
 (23)

Proof sketch. The path of DQ-GD and that of GD are related as (see Fig. 2, Lemma B.1)

$$\hat{\boldsymbol{x}}_t = \boldsymbol{x}_t - \eta \boldsymbol{e}_{t-1} \tag{24}$$

Comparing (24) and Line 4 in Algorithm 1, we see that  $z_t = x_t$ , i.e., DQ-GD computes the gradient at the unquantized trajectory  $\{x_t\}$ . The convergence guarantee of GD [67, Th. 1.4], [66, Theorem 2.1.15] (Lemma B.2) controls the difference between the first term in the recursion (24) and the optimizer  $x_f^*$ . To bound the second term in (24), we observe using (19) that for any  $r_t > 0$  in (7),

$$\max_{\boldsymbol{u} \in \mathcal{B}(r_t)} \| \mathsf{q}_t(\boldsymbol{u}) - \boldsymbol{u} \| = r_t \max_{\boldsymbol{u} \in \mathcal{B}(1)} \| \mathsf{q}(\boldsymbol{u}) - \boldsymbol{u} \|$$
(25)
$$= r_t \rho_n 2^{-R},$$
(26)

$$=r_t \rho_n 2^{-R}, \tag{26}$$

i.e. quantizer  $q_t$  used at iteration t has dynamic range  $r_t$ and covering radius (26). To complete the proof, we show by induction that with  $r_t$  in (22), the input  $u_t$  to the quantizer  $q_t$  generated by Algorithm 1 always lies within  $\mathcal{B}(r_t)$ . Since recurrence relation (22) represents a geometric sequence, (26) implies that the quantization error decays exponentially fast. The stepsize (20) is optimal both for GD [66, Theorem 2.1.15] and for DO-GD. See Appendix B-A for details.

The bound in (23) exhibits a phase-transition behavior: at any  $R \geq \log_2 \frac{\rho_n}{\sigma_{\rm GD}}$ , achieving the contraction factor of unquantized GD is possible, while at any  $R < \log_2 \frac{\rho_n}{\sigma_{\rm GD}}$ , the achievable contraction factor is only  $\rho_n 2^{-R} = \frac{\mathsf{d}(\mathsf{q})}{r(\mathsf{q})}$ . The algorithm converges linearly as long as  $\rho_n 2^{-R} < 1$ .

A common approach to quantizing descent algorithms [15]— [18], [55], [69] we refer to as naive quantization has the worker directly quantize the gradient of its current iterate. Applied to GD, it leads to the Naively Quantized Gradient Descent (NQ-GD) with the quantizer input (cf. (9))

$$\boldsymbol{u}_t = \nabla f(\hat{\boldsymbol{x}}_t). \tag{27}$$

In Theorem V.1 in Section V below, we show that

$$\sigma_{\text{NQ-GD}}(n,R) \le \sigma_{\text{GD}} + \frac{2\kappa}{\kappa+1} \rho_n 2^{-R}.$$
 (28)

which is strictly greater than (23).

In Fig. 3, we numerically compare the contraction factor of DQ-GD (Algorithm 1), the NQ-GD, and the unquantized GD (8) on least-squares problems

$$f(\boldsymbol{x}) = \frac{1}{2} \|\boldsymbol{y} - \mathbf{A}\boldsymbol{x}\|^2$$
 (29)

where  $y \in \mathbb{R}^m, \mathbf{A} \in \mathbb{R}^{m \times n}$ , with  $m \geq n$ . We generate 500 matrices A's with i.i.d. standard normal entries, one for each y, and rescale the spectrum of A so that it has a prescribed condition number  $\kappa$ . We also run the algorithm on the real-world least-squares matrix ash331 extracted from the online repository *SuiteSpare* [74]. For each per-dimension quantization rate  $R \geq 1$ , we generate 500 instances of the vector  $\boldsymbol{y}$  and  $\hat{\boldsymbol{x}}_0$  with i.i.d. standard normal entries. We run the

<sup>&</sup>lt;sup>3</sup>Covering efficiency introduced in (19) extends the notion of covering efficiency of an infinite lattice [73], which measures how well that lattice covers the whole space, to bounded-domain quantizers.

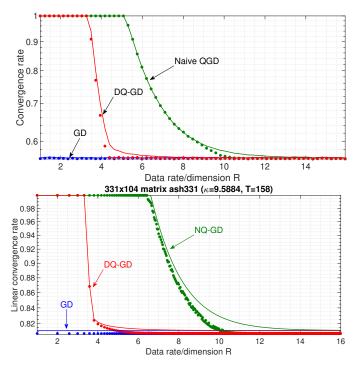


Fig. 3: Empirical contraction factors (as circles) and corresponding upper bounds (21), (23), and (28) (as lines).

iterative algorithms for as many iterations T as possible until reaching the machine's floating point precision, and report the average contraction factor. We use the uniform scalar quantizer for the ease of implementation and take as a consequence a space-filling loss of  $\sqrt{n}$ . For smaller values of the data rate R, quantized GD may not even converge as  $\sqrt{n}2^{-R} \geq 1$ . In that case, we clip off the contraction factor at 1 in the plots. We set the stepsize and the quantizer's dynamic range in the DQ-GD algorithm as prescribed by Theorem III.1, and in the NQ-GD algorithm as prescribed by Theorem V.1 in Section V below.

We observe that DQ-GD has a significantly faster contraction factor than NQ-GD, and that the empirical results closely track our analytical convergence bounds (23) and (28). The contraction factor of unquantized GD serves as a lower bound to both quantized algorithms.

Applying the error feedback mechanism of [10], [22], [44], developed for SGD, to GD results in an algorithm that forms the quantizer input as

$$\boldsymbol{u}_t = \nabla f(\hat{\boldsymbol{x}}_t) - \boldsymbol{e}_{t-1}. \tag{30}$$

Unlike DQ-GD (9), error feedback in (30) results in computing the gradient along the quantized trajectory  $\{\hat{x}_t\}$ , and it is unclear whether it can even improve upon NQ-GD (28) in the setting of our paper - nonstochastic GD with a worst-case performance criterion and without further assumptions on the quantizer (Appendix E).

#### C. DQ-AGD: convergence analysis

Unquantized accelerated gradient descent with stepsize

$$\eta = \eta_{\text{AGD}} \triangleq \frac{1}{L}$$
(31)

and momentum coefficient

$$\gamma = \gamma_{\text{AGD}} \triangleq \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \tag{32}$$

achieves contraction factor

$$\sigma_{\rm AGD} \triangleq \sqrt{1 - \frac{1}{\sqrt{\kappa}}}$$
 (33)

over  $\mathcal{F}_n$  (1) [33, Th. 3.18] (Lemma B.5), which improves the contraction factor of gradient descent  $\sigma_{\rm GD}=1-\frac{1}{\kappa}+O\left(\frac{1}{\kappa^2}\right)$  (21) to  $\sigma_{\rm AGD}=1-\frac{1}{2\sqrt{\kappa}}+O\left(\frac{1}{\kappa}\right)$ , a significant improvement if  $\kappa$  is large and optimal order-wise (the converse to the optimal contraction factor expands as  $1-\frac{4}{\sqrt{\kappa}}+O\left(\frac{1}{\kappa}\right)$  [66] (Lemma C.2) and is attained in  $\mathcal{F}_n^2$  by the heavy ball method [67] (Lemma B.8).

Denote for brevity the constant

$$\lambda \triangleq \left(1 + \gamma_{\text{AGD}} + \gamma_{\text{AGD}} \sigma_{\text{AGD}}^{-1}\right) \sqrt{\kappa + 1}.$$
 (34)

The following result extends (33) to DQ-AGD.

**Theorem III.2** (Convergence of DQ-AGD). Fix a dimensionn, rate-R quantizer q with dynamic range 1 and covering efficiency  $\rho_n$ . Then, Algorithm 2 with stepsize (31), momentum coefficient (32), and dynamic ranges  $r_{-2} = r_{-1} = 0$ ,

$$r_t = \sigma_{AGD}^t LD\lambda + (r_{t-1} + \gamma_{AGD}(r_{t-1} + r_{t-2})) \rho_n 2^{-R},$$
 (35)

t = 1, 2, ... in the definition of  $q_t$  (7) achieves the following contraction factor over  $\mathcal{F}_n$  (1):

$$\sigma_{\text{DQ-AGD}}(n, R) \le \max \left\{ \sigma_{\text{AGD}}, \rho_n 2^{-R} \phi(n, R, \gamma_{\text{AGD}}) \right\}$$
 (36)

where

$$\phi(n, R, \gamma) \triangleq \frac{1}{2}(1+\gamma) + \frac{1}{2}\sqrt{(1+\gamma)^2 + \frac{4\gamma}{\rho_n 2^{-R}}}.$$
 (37)

*Proof sketch.* The proof follows the roadmap of the proof of Theorem III.1 with the following complication. Where in Algorithm 1 the quantizer input depends on the previous quantization error  $e_{t-1}$ , the quantizer input in Algorithm 2 depends on the past two quantization errors  $e_{t-1}$  and  $e_{t-2}$  (Line 5). The resulting recursion (35) is a second-order linear non-homogeneous recurrence relation, which unlike (22) does not simply represent a geometric sequence. The characteristic polynomial of (35) is

$$p(r) \triangleq r^2 - r\rho_n 2^{-R} (1 + \gamma_{AGD}) - \rho_n 2^{-R} \gamma_{AGD}, \quad (38)$$

and  $\rho_n 2^{-R} \phi_{\mathrm{DQ-AGD}}(n,R)$  in (36) is its positive, larger-magnitude root. This implies that the quantization error decays with the contraction factor in the right side of (36). See Appendix B-B for details.

Define the functions

$$R_1(n,\gamma) \triangleq \log_2(1+2\gamma) + \log_2 \rho_n \tag{39}$$

$$R_2(n, \sigma, \gamma) \triangleq \log_2 \frac{(1+\gamma)\sigma + \gamma}{\sigma^2} + \log_2 \rho_n$$
 (40)

The achievability bound (36) exhibits two phase transitions. The first one is at  $\rho_n 2^{-R} \phi(n,R,\gamma_{\rm AGD}) < 1$ , which is equivalent to p(1) > 0: if

$$R > R_1(n, \gamma_{AGD})$$
 bits / dimension, (41)

then DQ-AGD enjoys linear convergence. The second one is at  $\rho_n 2^{-R} \phi(n, R, \gamma_{\rm AGD}) \leq \sigma_{\rm AGD}$ , which is equivalent to  $p(\sigma_{\rm AGD}) \geq 0$ : if

$$R \ge R_2(n, \sigma_{AGD}, \gamma_{AGD})$$
 bits / dimension, (42)

then there is no loss in the long-term convergence behavior of the DQ-AGD compared to AGD.

Curiously,  $R_1(n,0)$  and  $R_2(n,\sigma_{\rm GD},0)$  express the two phase transitions of DQ-DG that were determined in Section III-B.

D. DQ-HB: convergence analysis and numerical comparison Unquantized heavy ball method with stepsize

$$\eta = \eta_{\rm HB} \triangleq \left(\frac{2}{\sqrt{L} + \sqrt{\mu}}\right)^2$$
(43)

and momentum coefficient

$$\gamma = \gamma_{\rm HB} \triangleq \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}\right)^2 \tag{44}$$

achieves contraction factor

$$\sigma_{\rm HB} \triangleq \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \tag{45}$$

over  $\mathcal{F}_n^2$  (1) [67] (Lemma B.8), which is optimal among all gradient methods [66, Th. 2.1.13] (Lemma C.2).

The following convergence analysis of DQ-HB applies to smooth and strongly convex functions that are in addition twice continuously differentiable.

**Theorem III.3** (Convergence of DQ-HB). Fix a dimensionn, rate-R quantizer q with dynamic range 1 and covering efficiency  $\rho_n$ . Then, there exists a constant  $\alpha > 0$  such that Algorithm 3 with stepsize (43), momentum coefficient (44) and dynamic ranges  $r_{-1} = r_{-2} = 0$ ,

$$r_{t} = \sigma_{\text{HB}}^{t} t^{\alpha} e^{\alpha} \sqrt{2} LD + (r_{t-1} + \gamma_{\text{HB}} (r_{t-1} + r_{t-2})) \rho_{n} 2^{-R},$$
(46)

 $t = 1, 2, \dots$  achieves the following contraction factor over  $\mathcal{F}_n^2$ :

$$\sigma_{\rm DQ-HB}(n,R) \le \max\left\{\sigma_{\rm HB}, \rho_n 2^{-R}\phi(n,R,\gamma_{\rm HB})\right\},$$
 (47)  
where  $\phi(n,R,\gamma)$  is defined in (37).

*Proof sketch.* The proof is similar to the proof of Theorem III.2. The recurrence relation (46) differs from (35) in only the presence of the subexponential factor  $t^{\alpha}$ , which does not matter when we take  $t \to \infty$  to obtain (47). This factor arises from our nonasymptotic sharpening of Polyak's convergence result for the unquantized HB (Lemma B.8). See Appendix B-C for details.

DQ-HB exhibits two phase transitions, a behavior similar to DQ-HB and DQ-AGD. The two thresholds are given by  $R_1$  (39) and  $R_2$  (40) evaluated with  $\sigma = \sigma_{\rm HB}$  and  $\gamma = \gamma_{\rm HB}$ .

Plugging the parameters  $\gamma$  and  $\sigma$  into (40), we can infer that DQ-HB always has the largest  $R_2$  (40) for any condition number  $\kappa \geq 1$  among the three DQ schemes. On the other

hand, whether  $R_2$  of DQ-AGD is smaller than that of DQ-GD depends on whether  $\kappa$  is smaller than a threshold that is roughly 2.18. For the unquantized algorithms, contraction factor  $\sigma_{\rm HB}$  of HB is always the smallest among the three for any condition number  $\kappa \geq 1$ . On the other hand, whether  $\sigma_{\rm AGD}$  of AGD is smaller than  $\sigma_{\rm GD}$  of GD depends on whether  $\kappa$  is greater than a threshold that is roughly 11.83. For the differentially quantized algorithms, DQ-GD actually has the best convergence behavior in the transient regime where  $R > \log_2 \rho_n$  so that GD converges linearly and R is small enough so that DQ-HB does not yet outperform GD, i.e.,  $\rho_n 2^{-R} \phi(n, R, \gamma_{\rm HB}) > \sigma_{\rm GD}$ . This is because DQ-GD is the first among the three DQ algorithms to pass  $R_1$  (39) above which it has linear convergence.

In Fig. 4, we compare the performance of the differentially quantized algorithms DQ-HB (Algorithm 3), DQ-AGD (Algorithm 2) and DQ-GD (Algorithm 1) on least-squares problems (29). We use the same experimental setup as in Fig. 3, with the uniform scalar quantizer. The stepsize, the interpolation coefficient and the dynamic ranges are set to the values prescribed by Theorems III.1 (DQ-GD), III.2 (ADQ-GD), and III.3 (ADQ-HB). We set  $\alpha = 0$  for Algorithm 3, and DQ-HB still converges empirically for this parameter. We also record the performance of the corresponding unquantized gradient methods HB, AGD and GD. The curves exhibit the two phase transitions and comparative performance as discussed above. The level lines that the contraction factors of these DQ schemes rest on for  $R \geq R_2$  are almost the same as the corresponding linear convergence rates  $\sigma$  of their unquantized counterparts. We observe that there is a gap between the worst-case linear convergence rate  $\sigma_{AGD}$  that we design DQ-AGD to follow and the empirical convergence rate of AGD. This is because AGD applies for functions that are not necessarily twice continuously differentiable, and the leastsquares problems (29) happen not to be a worst-case problem class for AGD.

# IV. CONVERSES

# A. Quantized gradient descent algorithms

In this section, we characterize the optimal contraction factor achievable within class  $\mathcal{A}_{\mathrm{GD}}$  of quantized gradient descent algorithms, formally defined next.

**Definition IV.1** (Class  $\mathcal{A}_{\mathrm{GD}}$  of quantized gradient descent algorithms). A quantized gradient descent algorithm  $\mathbf{A} \in \mathcal{A}_{\mathrm{GD}}$  consists of a central server and an end worker. The algorithm is initialized with a collection of quantizers q indexed by rate R such that  $d(\mathbf{q}) \to 0$  (17) as  $R \to \infty$  and a sequence of dynamic ranges  $\{r_t\}_{t=1}^{\infty}$ . The worker has access to the function f. At each iteration t, the server first sends  $\hat{x}_t$  to the worker noiselessly, starting from some  $\hat{x}_0 \in \mathbb{R}^n$ . The worker then determines its gradient-access point  $z_t$  and its quantizer input  $u_t$  under the structural constraints

$$z_t \in \hat{x}_t + \operatorname{span}\left\{e_0, \dots, e_{t-1}\right\} \tag{48}$$

$$\boldsymbol{u}_t \in \nabla f(\boldsymbol{z}_t) + \operatorname{span}\left\{\boldsymbol{e}_0, \dots, \boldsymbol{e}_{t-1}\right\},$$
 (49)

where  $e_i \triangleq q_i - u_i$ , i = 0, ..., t - 1 are the past quantization errors before iteration t, and + denotes Minkowski's sum.

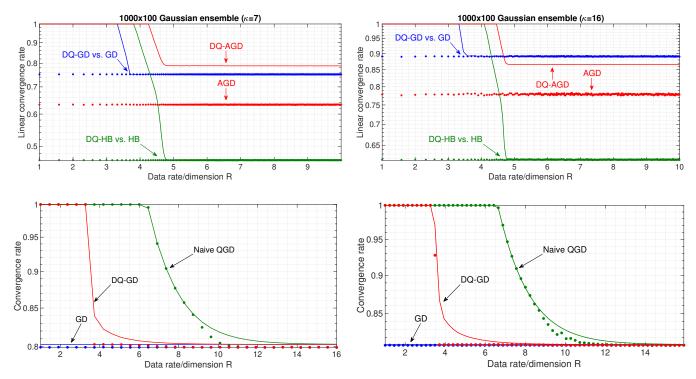


Fig. 4: Empirical contraction factors of various DQ algorithms (plotted as lines) and their corresponding unquantized counterparts (plotted as circles).

Upon receiving  $q_t = q_t(u_t)$  (7) from the worker, the server performs the update

$$\hat{\boldsymbol{x}}_{t+1} = \hat{\boldsymbol{x}}_t - \eta \boldsymbol{q}_t \tag{50}$$

with a fixed stepsize  $\eta > 0$ .

Due to conditions (48) and (49), if there is no quantization error at each iteration (i.e., if  $R=\infty$ ), then any quantized algorithm in  $\mathcal{A}_{\mathrm{GD}}$  reduces to the unquantized gradient descent. Both DQ-GD and NQ-GD fall in the class  $\mathcal{A}_{\mathrm{GD}}$ .

**Theorem IV.1** (Converse within class  $\mathcal{A}_{GD}$ ). The contraction factor achievable over functions  $f \in \mathcal{F}_n$  within class  $\mathcal{A}_{GD}$  of algorithms satisfies

$$\inf_{A \in \mathcal{A}_{GD}} \sigma_A(n, R) \ge \max \left\{ \sigma_{GD}, 2^{-R} \right\} \tag{51}$$

Proof sketch. We fix an  $A \in \mathcal{A}_{GD}$ , and we lower-bound the contraction factor it achieves at rate R in two different ways. On one hand, we show that A cannot converge faster than the unquantized GD. Then, we use an argument similar to [75] to craft a worst-case problem instance  $g \in \mathcal{F}_n$  for which the iterates of the unquantized GD satisfy  $||x_{t+1} - x_g^*|| = \sigma_{GD} ||x_t - x_g^*||$ , which ensures that  $\inf_{A \in \mathcal{A}_{GD}} \sigma_A(n, R) \geq \sigma_{GD}$  (Lemma C.1). On the other hand, we notice that if A is applied at dimension n and rate R, then the set  $\mathcal{S}_A \subseteq \mathbb{R}^n$  of all possible locations of the iterate  $\hat{x}_T$  after T iterations of A has cardinality at most  $2^{nRT}$ , and we apply a volume-division argument to claim that  $\inf_{A \in \mathcal{A}_{GD}} \sigma_A(n, R) \geq 2^R$ . See Appendix C-B for details.

Applying Theorem III.1 with Rogers-optimal quantizers with  $\rho_n \to 1$  [68, Th. 3] and juxtaposing with Theorem IV.1,

we characterize the optimal contraction factor achievable by quantized gradient descent in the limit of large problem dimension as

$$\lim_{n \to \infty} \inf_{A \in \mathcal{A}_{GD}} \sigma_A(n, R) = \max \left\{ \sigma_{GD}, 2^{-R} \right\}.$$
 (52)

In other words, DQ-DG achieves the best possible contraction factor within  $\mathcal{A}_{\mathrm{GD}}$ , in the limit of large problem dimension. This is rather remarkable: it means not only that DQ-DG compensates previous quantization errors optimally so that no rate is wasted, but that our convergence analysis in Theorem III.1 is tight enough to capture this optimality. Furthermore, notice that the right side of (52) is < 1 at any R > 0. This means that at any R > 0 however small, DQ-DG with Rogers-optimal quantizers converges linearly at a large enough problem dimension n, i.e. the first phase transition (39) dissappears.

Although the notion of a contraction factor (1) and thus the result in (52) are asymptotic in the number of iterations T, the achievability results in Appendix B-A and converse results in Appendix B-A used to derive (52) are nonasymptotic. They show that gap between the achievability and converse bounds on the finite-T counterpart of (52) is  $O\left(\frac{1}{T}\right)$ . Whether DQ-GD remains optimal at finite T remains an open problem.

#### B. Quantized gradient methods

All quantized algorithms considered in this paper fall in the following class.

**Definition IV.2** (Class  $\mathcal{A}_{GM}$  of quantized gradient methods). A quantized gradient method  $A \in \mathcal{A}_{GM}$  follows Definition IV.1 with (50) relaxed to

$$\hat{x}_{t+1} \in \hat{x}_0 + \text{span} \{q_0, \dots, q_t\}.$$
 (53)

In the absence of rate constraints, there are no quantization errors, i.e.  $e_t = \mathbf{0}$  for all t, and the class of quantized gradient methods reduces to the class of unquantized gradient methods satisfying

$$\boldsymbol{x}_{t+1} \in \boldsymbol{x}_0 + \operatorname{span} \left\{ \nabla f(\boldsymbol{x}_0), \dots, \nabla f(\boldsymbol{x}_t) \right\}.$$
 (54)

To present our converse result for  $A_{GM}$ , we consider functions f defined on the square-summable Hilbert space<sup>4</sup>

$$\mathbb{L}_2 \triangleq \left\{ \boldsymbol{x} = [\boldsymbol{x}(1), \boldsymbol{x}(2), \ldots] \colon \sum_{i=1}^{\infty} \boldsymbol{x}(i)^2 < \infty \right\}. \tag{55}$$

We say that continuously differentiable function  $f: \mathbb{L}_2 \to \mathbb{R}$  is in class  $\mathcal{F}_{\infty}$  if it is *L*-smooth,  $\mu$ -strongly convex and its minimizer is bounded, i.e., f satisfies i)-iii) in Section III-A.

To quantize an infinitely long vector  $u \in \mathbb{L}_2$  to  $q \in \mathbb{L}_2$ , we fix a free parameter  $n \in \mathbb{N}$ , apply a rate-R quantizer q in  $\mathbb{R}^n$  (6) to the first n coordinates of u, and set the remaining coordinates to 0, i.e.,

$$\begin{cases} [\boldsymbol{q}(1), \dots, \boldsymbol{q}(n)] &= \operatorname{q}\left([\boldsymbol{u}(1), \dots, \boldsymbol{u}(n)]\right) \\ \boldsymbol{q}(i) &= 0 \quad \forall i > n, \end{cases}$$
 (56)

where u(i) denotes *i*-th coordinate of vector  $u \in \mathbb{L}_2$ .

Although only n coordinates  $u \in \ell_2$  are quantized, we can still control the overall quantization error since in  $\mathbb{L}_2$ ,

$$\sum_{i>n} \boldsymbol{u}(i)^2 = o_n(1) \tag{57}$$

due to the Cauchy convergence criterion. Here  $o_n(1)$  denotes a function that vanishes as  $n\to\infty$ . Thus, (26) continues to hold for quantization in  $\mathbb{L}_2$  with  $\rho_n$  replaced by  $\rho_n+o_n(1)$ . It follows that the achievability bounds in Theorems III.2 and III.3 with with  $\rho_n$  replaced by  $\rho_n+o_n(1)$  apply to functions  $\mathbf{f}\in\mathcal{F}_\infty$ .

Contraction factor  $\sigma_A(n, R)$  over  $\mathcal{F}_{\infty}$  is defined in the same way as that over  $\mathcal{F}_n$  (1) except that n is now a parameter of the employed quantizer (like  $\rho_n$ ) rather than the dimension of the problem, and the total number of bits sent per iteration is nR, where R is the quantizer's rate (6).

**Theorem IV.2** (Converse within class  $\mathcal{A}_{GM}$ ). The contraction factor achievable over functions  $f \in \mathcal{F}_{\infty}$  within class  $\mathcal{A}_{GM}$  of algorithms satisfies

$$\inf_{A \in \mathcal{A}_{GM}} \sigma_A(n, R) \ge \max \left\{ \sigma_{HB}, 2^{-R} \right\}$$
 (58)

where  $\sigma_{\rm HB}$  is given in (45).

*Proof sketch.* The proof is similar to that of Theorem IV.1: we apply a volume-division argument to recover the  $2^{-R}$  in the

<sup>4</sup>We do so to take advantage of the sharpest converse in the literature on the convergence of unquantized gradient methods (54) [66] (Lemma C.2), which applies to functions on  $\mathbb{L}_2$ . Convergence lower bounds for smooth and strongly convex functions on  $\mathbb{R}^n$  (rather than  $\mathbb{L}_2$ ) are also known [76]. However, [76] considers only quadratic functions as objectives, and the considered class of iterative algorithms is more restrictive than (54) in that the next iterate  $\boldsymbol{x}_{t+1}$  depends on the past p terms  $\boldsymbol{x}_t, \ldots, \boldsymbol{x}_{t-p+1}$  for some fixed  $p \in \mathbb{N}$ .

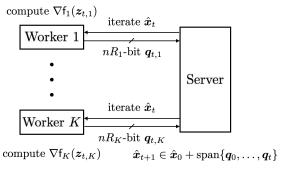


Fig. 5: K-worker quantized gradient method. At each iteration t, the server broadcasts the current iterate  $\hat{x}_t$ . Worker k computes the gradient at some point  $z_{t,k}$  that is a function of (but not necessarily equal to)  $\hat{x}_t$ . Then, worker k forms a descent direction  $q_{t,k}$  and pushes it back to the server under an  $nR_k$ -bit rate constraint.

right side of (58), and we apply a known result on unquantized gradient methods that states that the best contraction factor achievable over  $\mathcal{F}_{\infty}$  is that of the heavy ball method,  $\sigma_{\rm HB}$  [66] (Lemma C.2). See Appendix C-C.

Together, Theorems IV.2 and III.3 imply that for any  $R \geq R_2(\infty, \sigma_{\rm HB}, \gamma_{\rm HB})$ , DQ-HB attains the optimal contraction factor within  $\mathcal{A}_{\rm GM}$  (under the additional assumption that  $\mathsf{f} \in \mathcal{F}_\infty$  is twice continuously differentiable). The nonasymptotic achievability and converse results in Appendices B-C and C.2 used to show this result leave open the question of whether DQ-HB is optimal at finite number of iterations T, as they determine the finite-length analog of the optimal contraction factor only with accuracy  $O\left(\frac{\log T}{T}\right)$ .

# V. MULTI-WORKER GRADIENT METHODS

#### A. Problem setup

In empirical risk minimization [77], [78], the sample average of the loss function on the data points

$$f(\boldsymbol{x}) = \frac{1}{K} \sum_{k=1}^{K} f_k(\boldsymbol{x})$$
 (59)

arises as a substitute for the expected loss on the true data distribution that is often unknown. In multi-worker distributed empirical risk minimization, each worker has access to only one of the summands in (59), and they communicate to the parameter server under rate constraints. See Figure 5.

## B. Converses

Definitions IV.1 and IV.2 extend naturally to the K-worker setting. Converses in Theorems IV.1 and IV.2 extend verbatim to the multi-worker setting where the workers' rates satisfy the sum-rate constraint

$$\sum_{k=1}^{K} R_k \le R. \tag{60}$$

## C. Differential quantization

Differential quantization does not apply to K-worker quantized gradient methods since each worker does not know the local quantization errors stored by the others, and thus cannot guide the descent trajectory back to the unquantized path. Thus, whether (51) and (58) are attainable in the multiworker setting, and how each worker should optimally compensate its own past quantization errors, remain open problems.

## D. Naively Quantized Gradient Descent

The Naively Quantized Gradient Descent applies a common method of quantizing distributed gradient algorithms [15]–[18], [55], [69] in which each worker quantizes the gradient of the current iterate, to GD. It is summarized as Algorithm 4.

## **Algorithm 4:** K-worker NQ-GD

Our convergence result for NQ-GD holds under the following assumptions. We assume that continuously differentiable summands  $f_k$  in (59) are (i)  $L_k$ -smooth and (ii)  $\mu_k$ -strongly convex, and we continue to assume that (iii) the optimizer of f is bounded as in (15). Note that f is L-smooth and  $\mu$ -strongly convex with

$$L \triangleq \frac{1}{K} \sum_{k=1}^{K} L_k \quad \text{and} \quad \mu \triangleq \frac{1}{K} \sum_{k=1}^{K} \mu_k.$$
 (61)

Further, we focus on the *interpolation setting* [79]–[81] that assumes (iv)

$$\boldsymbol{x}_{\mathsf{f}}^* = \boldsymbol{x}_{\mathsf{f}_k}^* \quad \forall k = 1, \dots, K, \tag{62}$$

where

$$\boldsymbol{x}_{\mathsf{f}_k}^* \triangleq \operatorname*{arg\,min}_{\boldsymbol{x} \in \mathbb{R}^n} \mathsf{f}_k(\boldsymbol{x}).$$
 (63)

The interpolation setting is motivated by the observation that almost all local minima are also global in an over-parametrized neural network with a very large data dimension n [81], and is implied We denote by  $\mathcal{G}_n$  the class of functions f (59) that satisfy the assumptions (i)-(iv).

The minimum contraction factor achievable by K-worker NQ-GD under the sum rate constraint (60) is given by

$$\sigma_{\text{NO-GD}}(n,R)$$

$$\triangleq \inf_{\sum_{k=1}^{K} R_k \le R } \sup_{\mathbf{f} \in \mathcal{G}_n} \limsup_{T \to \infty} \left\| \hat{\boldsymbol{x}}_T(\{R_k\}_{k=1}^K) - \boldsymbol{x}_{\mathbf{f}}^* \right\|^{\frac{1}{T}}, \quad (64)$$

where  $\hat{x}_0(\{R_k\}_{k=1}^K), \hat{x}_1(\{R_k\}_{k=1}^K), \hat{x}_2(\{R_k\}_{k=1}^K), \dots$  is the sequence of iterates generated by NQ-GD (Algorithm 4) in response to  $f \in \mathcal{G}_n$  when the k-th worker operates at  $R_k$  bits per problem dimension,  $k = 1, \dots, K$ .

**Theorem V.1** (Convergence of K-worker NQ-GD). Fix a dimension-n, rate- $R_k$  quantizer  $q_k$  with dynamic range 1 and covering efficiency  $\rho_n$ . Set up the quantizer to be used by worker k at iteration t as

$$q_{t,k}(\cdot) = r_{t,k}q_k(\cdot/r_{t,k}), \tag{65}$$

where the dynamic ranges are given by

$$r_{t,k} = \left(\sigma_{\text{GD}} + \frac{\eta_{\text{GD}} \rho_n}{K} \sum_{k=1}^K \min\left\{\nu, L_k\right\}\right)^t L_k D, \quad (66)$$

and the optimum rate allocation is given by the waterfilling solution

$$R_k = |\log_2(L_k/\nu)|_+$$
 bits, (67)

where  $\nu$  is the water level found from the sum rate constraint

$$\sum_{k=1}^{K} |\log_2(L_k/\nu)|_{+} = R, \tag{68}$$

and  $|\cdot|_+ \triangleq \max\{0,\cdot\}$ . Then, Algorithm 4 with stepsize (20) achieves the following contraction factor over  $\mathcal{G}_n$  (64):

$$\sigma_{\text{NQ-GD}}(n, R) \le \sigma_{\text{GD}} + \frac{\eta_{\text{GD}} \rho_n}{K} \sum_{k=1}^{K} \min \{ \nu, L_k \}.$$
 (69)

*Proof.* Appendix 
$$D$$
.

According to (67), higher rates are allocated to users whose function gradients have higher Lipschitz constants, and if the Lipschitz constant is low enough in comparison to others no rate would be allocated at all. In the special case  $L_k \equiv L$ , (69) reduces to

$$\sigma_{\text{NQ-GD}}(n,R) \le \sigma_{\text{GD}} + \frac{2\kappa}{\kappa + 1} \frac{\rho_n}{2^{R/K}}.$$
 (70)

The bound in Theorem V.1 approaches the converse only in the limit of large R.

Without assumption (62) that all summands share the minimizer, NQ-GD converges only to a neighborhood of  $x_f^*$ ; the size of the neighborhood is controlled by the quantization error and vanishes as  $R \to \infty$  (Theorem D.2 in Appendix D).

# VI. CONCLUSION

This paper formalizes the problem of finding the optimal contraction factor achievable within a class of rate-constrained iterative optimization algorithms ((1), Definitions IV.1 and IV.2). We show information-theoretic converses to that fundamental limit (Theorem IV.1, Theorem IV.2).

We introduce the principle of differential quantization that posits that the quantizer's input shall be constructed in such a way as to guide the quantized algorithm's trajectory towards the unquantized trajectory. Applied to gradient descent (Algorithm 1), differential quantization leads to the contraction factor that is optimal within the class of quantized gradient descent algorithms ((52)). Thus, differential quantization leverages the memory of past quantized inputs in an optimal way, removing the impact of past quantization errors.

Beyond gradient descent, we apply differential quantization to gradient methods with momentum - the accelerated gradient descent (Algorithm 2) and the heavy ball method (Algorithm 3). In all three cases, differentially quantized algorithms attain the contraction factor of their unquantized counterparts as long as the data rate exceeds the corresponding threshold  $R_2$  (40).

Incidentally, in the course of the analysis, we provide a sharper bound on the convergence of the unquantized HB algorithm than available in the literature (Lemma B.8). We also provide a slightly more general worst-case problem instance for the unquantized GD than available (Lemma C.1).

Quantizers employed at each step have the same geometry (covering efficiency, (19)) but different resolution (covering radius, (18)). The resolution is controlled by scaling the quantizer's dynamic range (7). To attain the contraction factors in Theorems III.1, III.2, and III.3, the dynamic range is set to follow a recursion ((22), (35), (46)). That recursion shrinks the dynamic range at the fastest possible rate that still guarantees that the quantizer's input at each iteration falls within its dynamic range. This maximizes the usefulness of the bits exchanged at each iteration. While that recursion for DQ-GD (22) is simply a geometric sequence, those for DQ-AGD (35) and DQ-HB (46) are second-order linear non-homogeneous recurrence relations.

While DQ-GD attains the optimal contraction factor among quantized gradient descent algorithms (Definition IV.1) and DQ-HB attains the optimal contraction factor among all gradient methods, even unquantized, if  $R \geq R_2(n, \sigma_{\rm HB}, \gamma_{\rm HB})$  (40), it remains an open problem whether the contraction factor of  $2^{-R}$  dictated by the converse (Theorem IV.2) is achievable in the regime  $R_2(n, \sigma_{\rm GD}, 0) < R < R_2(n, \sigma_{\rm HB}, \gamma_{\rm HB})$  in the class of quantized gradient methods (Definition IV.2).

For multi-worker gradient descent, we provide a convergence result on naive quantization, in which the workers directly quantize their gradients, and show a waterfilling solution to optimize the allocation of data rates among the workers under the sum rate constraint (Theorem V.1). That result approaches the converse (Theorem IV.1) only in the limit  $R \to \infty$ , leaving open a tighter characterization of the optimum convergence factor in that scenario. Differential quantization does not directly apply to multi-worker optimization since the workers cannot compute the unquantized path without the knowledge of the local quantization errors stored by the others. We leave as future work the question of how they should optimally compensate their own quantization errors.

# ACKNOWLEDGMENT

The authors would like to thank Dr. Himanshu Tyagi for pointing out related works [19], [55]; Dr. Vincent Tan for bringing a known result on the worst-case contraction factor of unquantized GD [75] to our attention; Dr. Victor Kozyakin for a helpful discussion about joint spectral radius; and two anonymous reviewers for detailed comments.

# APPENDIX A DQ-GD WITH VARYING STEPSIZE

See Algorithm 5, below.

# Algorithm 5: DQ-GD with varying stepsize

```
1 Initialize e_{-1} = \hat{x}_0 = \mathbf{0}
2 for t = 0, 1, 2, \dots do
3 | Worker:
4 | \mathbf{z}_t = \hat{x}_t + \eta_{t-1}e_{t-1}
5 | \mathbf{u}_t = \nabla \mathbf{f}(\mathbf{z}_t) - (\eta_{t-1}/\eta_t)e_{t-1}
6 | \mathbf{q}_t = \mathbf{q}_t(\mathbf{u}_t)
7 | \mathbf{e}_t = \mathbf{q}_t - \mathbf{u}_t
8 | Server: \hat{x}_{t+1} = \hat{x}_t - \eta_t \mathbf{q}_t
9 end
```

#### APPENDIX B

# CONVERGENCE ANALYSES OF DQ ALGORITHMS

## A. Proof of Theorem III.1

As mentioned in the proof sketch, relation (24) is key to showing Theorem III.1. The next lemma, which applies to the more general version of the DQ-DG algorithm shown in Appendix A, establishes (24).

**Lemma B.1** (DQ-GD trajectory). Consider descent trajectories  $\{\hat{x}_t\}$  of Algorithm 5 and  $\{x_t\}$  of unquantized GD (8) with the same sequence of stepsizes  $\{\eta_t\}$  starting at the same location  $\hat{x}_0 = x_0$ . Then, at each iteration t,

$$\hat{x}_t = x_t - \eta_{t-1} e_{t-1}. \tag{71}$$

*Proof.* We prove (71) via mathematical induction.

- Base case: (71) holds for t = 0 since the starting location is the same.
- Inductive step: Suppose (71) holds for iteration t. First, the induction hypothesis, the quantizer input at Line 5 and the quantizer output at Line 6 of Algorithm 5 together imply

$$\boldsymbol{u}_t = \nabla f(\boldsymbol{x}_t) - \frac{\eta_{t-1}}{\eta_t} \boldsymbol{e}_{t-1}. \tag{72}$$

(We define  $0/0 \triangleq 0$  for the very first iteration when  $\eta_{-1} = 0$ .) We then have

$$\hat{\boldsymbol{x}}_{t+1} = \hat{\boldsymbol{x}}_t - \eta_t \boldsymbol{q}_t \tag{73}$$

$$=\hat{\boldsymbol{x}}_t - \eta_t(\boldsymbol{u}_t + \boldsymbol{e}_t) \tag{74}$$

$$= \hat{\boldsymbol{x}}_t - \eta_t \left( \nabla f(\boldsymbol{x}_t) - \frac{\eta_{t-1}}{\eta_t} \boldsymbol{e}_{t-1} \right) - \eta_t \boldsymbol{e}_t \quad (75)$$

$$= \left[ \boldsymbol{x}_t - \eta_t \nabla f(\boldsymbol{x}_t) \right] - \eta_t \boldsymbol{e}_t \tag{76}$$

$$= \boldsymbol{x}_{t+1} - \eta_t \boldsymbol{e}_t, \tag{77}$$

where (76) is due to the induction hypothesis.

Relation (71) implies for the constant stepsize  $\eta_t \equiv \eta$ 

$$\|\hat{\boldsymbol{x}}_t - \boldsymbol{x}_f^*\| \le \|\boldsymbol{x}_t - \boldsymbol{x}_f^*\| + \eta \|\boldsymbol{e}_{t-1}\|$$
 (78)

We use the contraction factor of unquantized GD to bound the first term of (78): **Lemma B.2** (Convergence of GD [66, Theorem 2.1.15]). For any L-smooth and  $\mu$ -strongly convex function f on  $\mathbb{R}^n$ , GD (8) with stepsize (20) satisfies

$$\|\boldsymbol{x}_t - \boldsymbol{x}_f^*\| \le \sigma_{\text{GD}}^t \|\boldsymbol{x}_0 - \boldsymbol{x}_f^*\|,$$
 (79)

where  $\sigma_{\rm GD}$  is defined in (21).

We use the following bound on quantization error to bound the second term in (78):

**Lemma B.3** (DQ-GD quantization error). Let  $f \in \mathcal{F}_n$ . Quantization errors  $\{e_t\}$  in Algorithm 1 with stepsize (20) and dynamic ranges (22) satisfy

$$\|\boldsymbol{e}_t\| \le r_t \, \rho_n 2^{-R} \tag{80}$$

$$\leq \max \left\{ \sigma_{\text{GD}}, \rho_n 2^{-R} \right\}^t LD \, \sigma_{\text{GD}} b_t,$$
 (81)

where

$$b_t \triangleq \begin{cases} \frac{\rho_n 2^{-R}}{|\sigma_{\text{GD}} - \rho_n 2^{-R}|} & \sigma_{\text{GD}} \neq \rho_n 2^{-R} \\ t + 1 & \sigma_{\text{GD}} = \rho_n 2^{-R} \end{cases}$$
(82)

*Proof.* Once we show that quantizer inputs  $\{u_t\}$  satisfy

$$\boldsymbol{u}_t \in \mathcal{B}(r_t), \tag{83}$$

(80) will follow from (26).

We prove (83) via induction.

• Base case: (83) holds for t = 0 since

$$\|\nabla f(\hat{x}_0) - e_{-1}\| = \|\nabla f(\hat{x}_0)\|$$
 (84)

$$< L \|\hat{\boldsymbol{x}}_0 - \boldsymbol{x}_{\mathsf{f}}^*\| \tag{85}$$

$$< LD,$$
 (86)

where (85) is due to  $\nabla f(x_f^*) = 0$  and L-smoothness (13), and (86) is due to assumption (15).

• Inductive step: Suppose (83) holds for iteration t. Applying triangle inequality and (71) to the expression in Line 5 yields

$$\|u_{t+1}\| < \|\nabla f(x_{t+1})\| + \|e_t\|.$$
 (87)

The first term is bounded as

$$\|\nabla f(x_{t+1})\| \le L \|x_{t+1} - x_f^*\|$$
 (88)

$$\leq \sigma_{\mathrm{GD}}^{t+1} L \| \boldsymbol{x}_0 - \boldsymbol{x}_{\mathsf{f}}^* \| \tag{89}$$

$$\leq \sigma_{\rm CD}^{t+1} LD,$$
 (90)

where (88) is due to  $\nabla f(x_f^*) = \mathbf{0}$  and L-smoothness (13); (89) is due to (79); and (90) is due to assumption (15). Quantization error term  $\|e_t\|$  in (87) is bounded by (26) due to the induction hypothesis. Plugging (90) and (26) into (87) gives

$$\|\boldsymbol{u}_{t+1}\| \le \sigma_{\text{GD}}^{t+1} L D + r_t \rho_n 2^{-R}$$
 (91)

$$=r_{t+1}, (92)$$

where (92) is due to the choice of the dynamic ranges (22).

This concludes the proof of (80). To establish (81), we unwrap recursion (22) as the geometric sum

$$r_{t} = LD \sum_{\tau=0}^{t} \sigma_{\text{GD}}^{\tau} \left(\rho_{n} 2^{-R}\right)^{t-\tau}$$

$$= LD \cdot \begin{cases} \sigma_{\text{GD}}^{t} \frac{1 - \left[\rho_{n} 2^{-R} / \sigma_{\text{GD}}\right]^{t+1}}{1 - \rho_{n} 2^{-R} / \sigma_{\text{GD}}} & \sigma_{\text{GD}} \neq \rho_{n} 2^{-R} \\ t + 1 & \sigma_{\text{GD}} = \rho_{n} 2^{-R} \end{cases}$$

$$\leq LD \cdot \begin{cases} \sigma_{\text{GD}}^{t} \left(1 - \frac{\rho_{n} 2^{-R}}{\sigma_{\text{GD}}}\right)^{-1} & \sigma_{\text{GD}} > \rho_{n} 2^{-R} \\ \left(\rho_{n} 2^{-R}\right)^{t} \left(\frac{\rho_{n} 2^{-R}}{\sigma_{\text{GD}}} - 1\right)^{-1} & \sigma_{\text{GD}} < \rho_{n} 2^{-R} \\ t + 1 & \sigma_{\text{GD}} = \rho_{n} 2^{-R} \end{cases}$$

$$(93)$$

where the bound for the case  $\sigma_{\mathrm{GD}} > \rho_n 2^{-R}$  is obtained by lower-bounding  $\left[\rho_n 2^{-R}/\sigma_{\mathrm{GD}}\right]^{t+1}$  by 0, and the bound for  $\sigma_{\mathrm{GD}} < \rho_n 2^{-R}$  is obtained by lower-bounding  $\left[\sigma_{\mathrm{GD}}/\rho_n 2^{-R}\right]^{t+1}$  by 0.

Putting together the results in Lemmas B.1, B.2, and B.3, we show the following nonasymptotic (in the iteration number) convergence result for the DQ-GD:

**Theorem B.1** (Convergence of DQ-GD). In the setting of Theorem III.1, the difference between the iterate and the optimizer at step t satisfies

$$\|\hat{\boldsymbol{x}}_t - \boldsymbol{x}_f^*\| \le \max \left\{ \sigma_{\text{GD}}, \, \rho_n 2^{-R} \right\}^t \left[ 1 + \eta_{\text{GD}} L b_{t-1} \right] D.$$
 (96)

*Proof.* Plugging (79) and (81) into (78), we obtain

$$\|\hat{\boldsymbol{x}}_t - \boldsymbol{x}_f^*\| \le \sigma_{\text{GD}}^t D$$

$$+ \max \left\{ \sigma_{\text{GD}}, \rho_n 2^{-R} \right\}^{t-1} \eta_{\text{GD}} LD \, \sigma_{\text{GD}} b_{t-1},$$
(97)

which leads to (96) by an elementary weakening.

Bound (23) in Theorem III.1 follows immediately by applying (96) to definition (1) of the contraction factor.

#### B. Proof of Theorem III.2

The proof follows steps similar to those in the proof of Theorem III.1. First, we prove that, as prescribed by the principle of differential quantization, DQ-AGD compensates quantization errors by directing the quantized trajectory back to the trajectory of AGD:

**Lemma B.4** (DQ-AGD trajectory). *Iterate sequences*  $\{\hat{y}_t, \hat{x}_t\}$  of Algorithm 2 and  $\{y_t, x_t\}$  of AGD (10) starting at the same location  $(\hat{y}_0, \hat{x}_0) = (y_0, x_0)$  are related as,  $\forall t = 0, 1, 2, ...$ ,

$$\hat{\mathbf{y}}_t = \mathbf{y}_t - \eta \mathbf{e}_{t-1} \tag{98}$$

$$\hat{x}_{t} = x_{t} - \eta e_{t-1} - \eta \gamma \left( e_{t-1} - e_{t-2} \right). \tag{99}$$

*Proof.* We prove (98) and (99) by induction.

• Base case: (98) and (99) hold for t=0 since both algorithms start at the same location.

• Inductive step: Line 9 yields

$$\hat{\boldsymbol{y}}_{t+1} = \hat{\boldsymbol{x}}_t - \eta \boldsymbol{q}_t \tag{100}$$

$$= x_t - \eta \left[ e_{t-1} + \gamma \left( e_{t-1} - e_{t-2} \right) \right]$$
 (101)

$$-\eta \left[\nabla \mathsf{f}(\boldsymbol{x}_{t}) - \left[\boldsymbol{e}_{t-1} + \gamma \left(\boldsymbol{e}_{t-1} - \boldsymbol{e}_{t-2}\right)\right] + \boldsymbol{e}_{t}\right]$$

$$= y_{t+1} - \eta e_t, \tag{102}$$

where (101) is due to induction hypothesis (99) and Lines 5 and 7, and (102) is due to (10). On the other hand, plugging (98) and (102) into Line 10 yields

$$\hat{x}_{t+1} = \hat{y}_{t+1} + \gamma \left( \hat{y}_{t+1} - \hat{y}_{t} \right)$$

$$= y_{t+1} + \gamma \left( y_{t+1} - y_{t} \right) - \eta \left[ e_{t} + \gamma \left( e_{t} - e_{t-1} \right) \right]$$

$$= x_{t+1} - \eta \left[ e_{t} + \gamma \left( e_{t} - e_{t-1} \right) \right],$$
(104)

where (104) is due to (11).

Via triangle inequality, relation (98) implies

$$\|\hat{y}_t - x_f^*\| \le \|y_t - x_f^*\| + \eta \|e_{t-1}\|$$
 (105)

The following simple corollary to a known bound on the contraction factor of unquantized AGD controls the first term of (105):

**Lemma B.5** (Convergence of AGD). For any L-smooth and  $\mu$ -strongly convex function f on  $\mathbb{R}^n$ , AGD ((10), (11) starting at  $\mathbf{x}_0 = \mathbf{y}_0$ ) with stepsize (31) and momentum coefficient (32) satisfies

$$\|\boldsymbol{y}_{t} - \boldsymbol{x}_{\mathsf{f}}^{*}\| \le \sigma_{\mathrm{AGD}}^{t} \sqrt{\kappa + 1} \|\boldsymbol{x}_{0} - \boldsymbol{x}_{\mathsf{f}}^{*}\| \tag{106}$$

$$\|\boldsymbol{x}_{t} - \boldsymbol{x}_{\mathsf{f}}^{*}\| \leq \sigma_{\mathrm{AGD}}^{t} \lambda \|\boldsymbol{x}_{0} - \boldsymbol{x}_{\mathsf{f}}^{*}\|. \tag{107}$$

*Proof.* According to [33, Theorem 3.18],

$$f(y_t) - f(x_f^*) \le \frac{L + \mu}{2} \sigma_{AGD}^{2t} ||x_0 - x_f^*||^2.$$
 (108)

Convergence bound (106) w.r.t.  $\{y_t\}$  is due to (108) and

$$f(x) - f(x_f^*) \ge \frac{\mu}{2} ||x - x_f^*||^2$$
 (109)

implied by  $\mu$ -strong convexity [66, Theorem 2.1.8]. On the other hand, applying triangle inequality to (11) yields

$$\|\boldsymbol{x}_{t+1} - \boldsymbol{x}_{f}^{*}\| < (1+\gamma) \|\boldsymbol{y}_{t+1} - \boldsymbol{x}_{f}^{*}\| + \gamma \|\boldsymbol{y}_{t} - \boldsymbol{x}_{f}^{*}\|,$$
 (110)

and (107) then follows from (106) and (110). 
$$\Box$$

The following bound on the quantization error controls the second term in (78):

**Lemma B.6** (DQ-AGD quantization error). Let  $f \in \mathcal{F}_n$ . Quantization errors  $\{e_t\}$  in Algorithm 2 with stepsize (31), momentum coefficient (32) and dynamic ranges (35) satisfy

$$\|\boldsymbol{e}_t\| \le r_t \rho_n 2^{-R} \tag{111}$$

$$\leq \sigma_{\text{AGD}}^t c_0 + \phi_+^t c_+ + \phi_-^t c_-,$$
 (112)

where  $\phi_{\pm} \triangleq \phi_{\pm}(\gamma_{AGD})$  with

$$\phi_{\pm}(\gamma) \triangleq \rho_n 2^{-R} \left( \frac{1}{2} (1+\gamma) \pm \frac{1}{2} \sqrt{(1+\gamma)^2 + \frac{4\gamma}{\rho_n 2^{-R}}} \right)$$
(113)

and  $c_0, c_+, c_-$  are specified below in (118), (122) and (123) respectively.

*Proof.* Like in the proof of Lemma B.3, to establish (111), in view of (26) it is enough to show (83). Applying triangle inequality and (99) to the expression in Line 5 yields

$$\|\mathbf{u}_t\| \le \|\nabla f(\mathbf{x}_t)\| + \|\mathbf{e}_{t-1}\| + \gamma (\|\mathbf{e}_{t-1}\| + \|\mathbf{e}_{t-2}\|).$$
 (114)

The first term in (114) is bounded as

$$\|\nabla f(\boldsymbol{x}_t)\| \le L \|\boldsymbol{x}_t - \boldsymbol{x}_f^*\| \tag{115}$$

$$\leq \sigma_{\text{AGD}}^t LD\lambda$$
 (116)

where (115) is due to L-smoothness (13), and (116) is due to (107). Plugging (116) into (114) and applying (26) to bound quantization error terms in (114), we conclude via induction (similar to the proof of Lemma B.3) that setting the sequence of dynamic ranges recursively as (35) ensures (83).

To show (112), we proceed to solve recursion (35). This step is significantly different from the corresponding step in the proof of Lemma B.3 since now  $r_t$  depends not just on  $r_{t-1}$  but also on  $r_{t-2}$ . More precisely, recursion (35) is a second-order linear non-homogeneous recurrence relation.

• Particular solution: Plugging the candidate

$$p_t = \sigma_{\text{AGD}}^t c_0 \tag{117}$$

into (35), we solve for the constant

$$c_0 = \frac{\sigma_{\text{AGD}}^2}{\mathsf{p}(\sigma_{\text{AGD}})} LD\lambda,\tag{118}$$

where p(r) is the characteristic polynomial in (38) associated with recursion (35).

• Homogeneous solution: Since  $\phi_+$  and  $\phi_-$  are roots of the quadratic polynomial in (38), the homogeneous solution is given by

$$\phi_t = \phi_+^t c_+ + \phi_-^t c_-, \tag{119}$$

• General solution: constants  $c_+,c_-$  in (119) are determined by plugging initial conditions  $r_{-2}=r_{-1}=0$  into the general solution

$$r_t = p_t + \phi_t \tag{120}$$

$$= \sigma_{\text{AGD}}^t c_0 + \phi_+^t c_+ + \phi_-^t c_-, \tag{121}$$

which are

$$c_{+} = -\frac{c_{0}\phi_{+}^{2}}{\sigma_{\text{AGD}}^{2}} \frac{\sigma_{\text{AGD}} - \phi_{-}}{\phi_{+} - \phi_{-}}$$
(122)

$$c_{-} = \frac{c_0 \phi_{-}^2}{\sigma_{\text{AGD}}^2} \frac{\sigma_{\text{AGD}} - \phi_{+}}{\phi_{+} - \phi_{-}}.$$
 (123)

Lemmas B.4, B.5, and B.6 lead to the following finite-*t* convergence bound for the DQ-AGD:

**Theorem B.2** (Convergence of DQ-AGD). In the setting of Theorem III.2, the difference between the iterate and the optimizer at step t satisfies

$$\|\hat{y}_t - x_f^*\| \le \sigma_{AGD}^t c + \phi_+^{t-1} c_+ \eta + \phi_-^{t-1} c_- \eta$$
 (124)

where  $c = \sqrt{\kappa + 1}D + \eta c_0 \sigma_{AGD}^{-1}$ , and  $c_0, c_+, c_-$  are specified in (118), (122) and (123) respectively.

*Proof.* Plugging (106) and (112) into (105) immediately leads to (124).  $\Box$ 

Note that if  $\sigma_{\rm AGD} \geq \phi_+$ , which is equivalent to  $R \geq R_2$  (40), then  $c_0 \geq 0$ ,  $c_+ \leq 0$ ,  $c_- \geq 0$ ; and  $c_0 \leq 0$ ,  $c_+ \geq 0$ ,  $c_- \leq 0$  otherwise. Asymptotic convergence guarantee (36) in Theorem III.2 follows by plugging (124) into definition (1) of the contraction factor.

#### C. Proof of Theorem III.3

The proof of the DQ-HB convergence result in Theorem III.3 follows the same recipe as the proofs of Theorems III.1 and III.2. Lemma B.7 below states that the path of DQ-HB tracks that of the unquantized HB.

**Lemma B.7** (DQ-HB trajectory). Path  $\{\hat{x}_t\}$  of DQ-HB (Line 7) and path  $\{x_t\}$  of HB (12) starting at the same location  $\hat{x}_{-1} = x_{-1} = \hat{x}_0 = x_0$  are related as,  $\forall t = 0, 1, 2, ...$ ,

$$\hat{\boldsymbol{x}}_t = \boldsymbol{x}_t - \eta \boldsymbol{e}_{t-1}. \tag{125}$$

*Proof.* We prove (125) via induction.

- Base case: (125) holds for t = 0 by the initialization  $e_{-2} = e_{-1} = 0$  in Line 1 of Algorithm 3 and since the starting location is the same.
- Inductive step: Plugging expressions on Line 5 and 7 into Line 8 yields

$$\hat{x}_{t+1} \qquad (126)$$

$$= \hat{x}_t - \eta \mathbf{q}_t + \gamma (\hat{x}_t - \hat{x}_{t-1}) \qquad (127)$$

$$= \mathbf{x}_t - \eta \mathbf{e}_{t-1} - \eta [\nabla f(\mathbf{x}_t) - [\mathbf{e}_{t-1} + \gamma (\mathbf{e}_{t-1} - \mathbf{e}_{t-2})] + \mathbf{e}_t] + \gamma (\mathbf{x}_t - \mathbf{x}_{t-1} - \eta (\mathbf{e}_{t-1} - \mathbf{e}_{t-2})) \qquad (128)$$

$$= \mathbf{x}_t - \eta \nabla f(\mathbf{x}_t) + \gamma (\mathbf{x}_t - \mathbf{x}_{t-1}) - \eta \mathbf{e}_t \qquad (129)$$

$$= \boldsymbol{x}_{t+1} - \eta \boldsymbol{e}_t, \tag{130}$$

where (128) is due to the induction hypothesis and (130) is due to (12).

Applying triangle inequality to (125) gives

$$\|\hat{x}_t - x_f^*\| \le \|x_t - x_f^*\| + \eta \|e_{t-1}\|.$$
 (131)

The convergence result for unquantized HB in Lemma B.8, below, controls the first term in the right side of (131). Lemma B.8 is a nonasymptotic refinement of Polyak's original convergence result [67, Th. 1, Sec. 3.2]. Unlike the original, it does not require that the algorithm starts "sufficiently close" to the minimizer (i.e., it establishes global rather than local convergence), and it also clarifies that the subexponential factor in the bound is polynomial in t (see (132), below). This refinement is made possible by a result on joint spectral radius due to Wirth [82] that is more recent than [67, Th. 1, Sec. 3.2].

**Lemma B.8** (Convergence of HB). For any L-smooth,  $\mu$ -strongly convex, twice continuously differentiable function f on

 $\mathbb{R}^n$ , there exists a constant  $\alpha > 0$  such that the HB algorithm (12) starting at  $\mathbf{x}_{-1} = \mathbf{x}_0$  with stepsize (43) and momentum coefficient (44) satisfies

$$\|\boldsymbol{x}_t - \boldsymbol{x}_f^*\| \le \sigma_{HB}^t t^{\alpha} e^{\alpha} \sqrt{2} \|\boldsymbol{x}_0 - \boldsymbol{x}_f^*\|.$$
 (132)

*Proof.* Iterative process (12) can be written in the form

$$\begin{bmatrix} \boldsymbol{x}_{t+1} - \boldsymbol{x}_{f}^{*} \\ \boldsymbol{x}_{t} - \boldsymbol{x}_{f}^{*} \end{bmatrix}$$

$$= \begin{bmatrix} \boldsymbol{x}_{t} + \gamma(\boldsymbol{x}_{t} - \boldsymbol{x}_{t-1}) - \boldsymbol{x}_{f}^{*} \\ \boldsymbol{x}_{t} - \boldsymbol{x}_{f}^{*} \end{bmatrix} - \eta \begin{bmatrix} \nabla f(\boldsymbol{x}_{t}) \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} (1+\gamma)\boldsymbol{I} & -\gamma\boldsymbol{I} \\ \boldsymbol{I} & 0 \end{bmatrix} \begin{bmatrix} \boldsymbol{x}_{t} - \boldsymbol{x}_{f}^{*} \\ \boldsymbol{x}_{t-1} - \boldsymbol{x}_{f}^{*} \end{bmatrix} - \eta \begin{bmatrix} \nabla^{2}f(\boldsymbol{v}_{t})(\boldsymbol{x}_{t} - \boldsymbol{x}_{f}^{*}) \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} (1+\gamma)\boldsymbol{I} & -\gamma\boldsymbol{I} - \eta\nabla^{2}f(\boldsymbol{v}_{t}) \\ \boldsymbol{I} & 0 \end{bmatrix} \begin{bmatrix} \boldsymbol{x}_{t} - \boldsymbol{x}_{f}^{*} \\ \boldsymbol{x}_{t-1} - \boldsymbol{x}_{f}^{*} \end{bmatrix},$$

$$(135)$$

where (134) holds for some  $v_t$  on the line segment between  $x_t$  and  $x_t^*$  by the mean value theorem, since f is twice continuously differentiable by the assumption. Denoting the matrix in (135) by  $A_t$ , we unroll the recursion as

$$\begin{bmatrix} \boldsymbol{x}_{t+1} - \boldsymbol{x}_{f}^{*} \\ \boldsymbol{x}_{t} - \boldsymbol{x}_{f}^{*} \end{bmatrix} = \boldsymbol{A}_{t} \cdot \boldsymbol{A}_{t-1} \cdot \ldots \cdot \boldsymbol{A}_{0} \begin{bmatrix} \boldsymbol{x}_{0} - \boldsymbol{x}_{f}^{*} \\ \boldsymbol{x}_{-1} - \boldsymbol{x}_{f}^{*} \end{bmatrix}. \quad (136)$$

It follows that

$$\left\| \begin{bmatrix} \boldsymbol{x}_{t+1} - \boldsymbol{x}_{\mathsf{f}}^* \\ \boldsymbol{x}_t - \boldsymbol{x}_{\mathsf{f}}^* \end{bmatrix} \right\|_2 \le \left\| \boldsymbol{A}_t \cdot \ldots \cdot \boldsymbol{A}_0 \right\|_2 \left\| \begin{bmatrix} \boldsymbol{x}_0 - \boldsymbol{x}_{\mathsf{f}}^* \\ \boldsymbol{x}_{-1} - \boldsymbol{x}_{\mathsf{f}}^* \end{bmatrix} \right\|_2. \quad (137)$$

It is shown in [82, Lemma 2.3] that if matrices  $A_1, \ldots, A_t$  all belong to a bounded set A, then there exists a constant  $\alpha > 0$  such that  $\forall t = 0, 1, \ldots$ ,

$$\|\boldsymbol{A}_t \cdot \dots \cdot \boldsymbol{A}_0\|_2 \le \rho(\mathcal{A})^{t+1} (t+1)^{\alpha} e^{\alpha}$$
 (138)

where

$$\rho(\mathcal{A}) \triangleq \limsup_{t \to \infty} \sup_{t} \rho(\mathbf{A}_t), \tag{139}$$

where  $\rho(A_t)$  is the spectral radius of  $A_t$ . It is shown in [67, Proof of Th. 1, Sec. 3.2] that if

$$\gamma = \max\left\{ \left(1 - \sqrt{\eta L}\right)^2, \ \left(1 - \sqrt{\eta \mu}\right)^2 \right\}, \tag{140}$$

then

$$\rho(\mathbf{A}_t) \le \sqrt{\gamma}.\tag{141}$$

With the optimal choice of  $\eta$  (43) the right side of (141) is equal to  $\sigma_{\rm HB}$  (45).

In our setting  $\mathcal A$  is bounded since for twice continuously differentiable functions, L-smoothness and  $\mu$ -strong convexity are equivalent to

$$\mu \mathbf{I} \prec \nabla^2 \mathsf{f}(\mathbf{v}) \prec L \mathbf{I},$$
 (142)

in the positive semidefinite order, thus (138) applies. Inequality (132) follows after applying (141) to (138) and the latter to (137).

Remark B.1. Polyak's convergence result [67, Th. 1, Sec. 3.2] guarantees only the existence of D > 0 such that for all starting points  $x_{-1}, x_0$  with  $||x_{\mathbf{f}}^* - x_{-1}|| \le D$ ,  $||x_{\mathbf{f}}^* - x_0|| \le D$ 

and all  $0 < \epsilon < 1 - \sigma_{\rm HB}$ , there exists a c > 0 such that (cf. (132))

$$\|\boldsymbol{x}_t - \boldsymbol{x}_f^*\| \le c(\sigma_{HB} + \epsilon)^t \tag{143}$$

under the same assumptions on f and the same stepsize and momentum coefficient as in Lemma B.8. This is a local convergence result because convergence is not guaranteed for any starting location but only for locations in a small enough neighbourhood of  $x_f^*$ . Furthermore, (132) refines the subexponential factor in (143).

We ensure that the quantization error term  $||e_{t-1}||$  in (131) decays exponentially fast by adjusting the sequence of dynamic ranges (7) carefully:

**Lemma B.9** (DQ-HB quantization error). Let  $f \in \mathcal{F}_n^2$ . Quantization errors  $\{e_t\}$  in Algorithm 3 with with stepsize (43), momentum coefficient (44) and dynamic ranges (46) satisfy

$$\|\boldsymbol{e}_t\| \le r_t \rho_n 2^{-R} \tag{144}$$

$$\leq (\sigma_{\rm HB}^t c_0 + \phi_+^t c_+ + \phi_-^t c_-) t^{\alpha},$$
 (145)

where  $\phi_{\pm} \triangleq \phi_{\pm}(\gamma_{\rm HB})$  (113), and  $c_0, c_+, c_-$  are as in (118), (122) and (123) respectively, with  $LD\lambda$  replaced by  $e^{\alpha}\sqrt{2}LD$ , and  $\sigma_{\rm AGD}$  by  $\sigma_{\rm HB}$ .

*Proof.* Applying triangle inequality and (125) to the expression in Line 5 yields the same expression as in the analysis of DQ-AGD (114):

$$\|\boldsymbol{u}_t\| \le \|\nabla f(\boldsymbol{x}_t)\| + \|\boldsymbol{e}_{t-1}\| + \gamma (\|\boldsymbol{e}_{t-1}\| + \|\boldsymbol{e}_{t-2}\|).$$
 (146)

The first term in (146) is bounded as

$$\|\nabla f(\boldsymbol{x}_t)\| \le L \|\boldsymbol{x}_t - \boldsymbol{x}_f^*\| \tag{147}$$

$$\leq \sigma_{\rm HB}^t \, t^\alpha e^\alpha \sqrt{2} \, LD \tag{148}$$

where (115) is due to L-smoothness (13), and (148) is due to (132). Applying the argument used to show (111) in the proof of Lemma B.6 leads to (144).

To show (145), consider the recursion  $r'_{-1} = r'_{-2} = 0$ ,

$$r'_{t} = \sigma_{\text{HB}}^{t} e^{\alpha} \sqrt{2} LD + \left( r'_{t-1} + \gamma_{\text{HB}} (r'_{t-1} + r_{t-2}) \right) \rho_{n} 2^{-R},$$
(149)

We show by strong induction that

$$r_t \le t^{\alpha} r_t',\tag{150}$$

where  $r_t$  solves (46). Base case  $r_0 = r_0'$  holds by the initial conditions. Assuming that (150) holds for  $1, \ldots, t-1$ , we establish (150) for t using (46) and the fact that  $t^{\alpha}$  is increasing in t:

$$r_{t} \leq \sigma_{\text{HB}}^{t} t^{\alpha} e^{\alpha} \sqrt{2} LD$$

$$+ \left( (t-1)^{\alpha} r'_{t-1} + \gamma_{\text{HB}} ((t-1)^{\alpha} r'_{t-1} + (t-2)^{\alpha} r'_{t-2}) \right) \rho_{n} 2^{-R}$$

$$< t^{\alpha} r'_{t}.$$
(152)

Using (144) and (150), we can show (145) by solving the recursion (149). But this is the same recursion as in (35), up to the constants, thus the solution in the proof of Lemma B.6 applies.

We now apply Lemmas B.7, B.8 and B.9 to state a finite-*t* refinement of Theorem III.3.

**Theorem B.3** (Convergence of DQ-HB). In the setting of Theorem III.3, the difference between the iterate and the optimizer at step t satisfies

$$\|\hat{\boldsymbol{x}}_t - \boldsymbol{x}_f^*\| \le \left(\sigma_{HB}^t c + \phi_{\perp}^{t-1} c_{\perp} + \phi_{\perp}^{t-1} c_{-}\right) t^{\alpha}$$
 (153)

where  $c = e^{\alpha}\sqrt{2}D + \eta c_0\sigma_{\rm HB}^{-1}$ , and  $\phi_+, \phi_-, c_0, c_+, c_-$  are as in Lemma B.9.

*Proof.* Plugging (132) and (145) into (131) and using  $(t-1)^{\alpha} < t^{\alpha}$  leads to (153).

# APPENDIX C CONVERSES

#### A. Converse for unquantized GD

As explained in the proof sketch, the lower bound in (51) is a combination of an unquantized GD converse and a volume-division converse. The former relies on the following converse result, obtained by constructing a least-square problem instance that satisfies Nesterov's upper bound in Lemma B.2 with equality.

**Lemma C.1** (Optimality of  $\sigma_{GD}$ ). Consider GD (8) with starting point  $x_0$  and any constant stepsize  $\eta$ . Then, there exists a problem instance  $f \in \mathcal{F}_n$  such that the distance to the optimizer at each iteration t of GD satisfies

$$\|\boldsymbol{x}_{t+1} - \boldsymbol{x}_{f}^*\| \ge \sigma_{GD} \|\boldsymbol{x}_{t} - \boldsymbol{x}_{f}^*\|,$$
 (154)

with equality if  $\eta$  is the optimal stepsize given in (20).

*Proof of Lemma C.1.* We first find an  $x_f^* \in \mathcal{B}(D)$  such that

$$\|\boldsymbol{x}_0 - \boldsymbol{x}_f^*\| > D \tag{155}$$

and then construct a least-squares problem instance  $f \in \mathcal{F}_n$  (29) that admits  $x_f^*$  as a unique minimizer and satisfies (154). Note that f is

$$\sigma_1^2(\mathbf{A})$$
-smooth and  $\sigma_n^2(\mathbf{A})$ -strongly convex (156)

where we denote by  $\sigma_i(\mathbf{A})$  the *i*-th largest singular value of matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . The gradient of f at iteration t is

$$\nabla f(\boldsymbol{x}_t) = \mathbf{A}^{\top} (\mathbf{A} \boldsymbol{x}_t - \boldsymbol{y}). \tag{157}$$

The first-order optimality condition  $abla f(oldsymbol{x}_{\mathsf{f}}^*) = \mathbf{0}$  implies

$$\mathbf{A}^{\top} y = \mathbf{A}^{\top} \mathbf{A} x_{\mathsf{f}}^{*}. \tag{158}$$

To each  $x_f^* \in \mathcal{B}(D)$  that satisfies (155), there corresponds a  $y \in \mathbb{R}^m$  such that (158) holds. This is because  $m \geq n$ , i.e. we have more degrees of freedom than the problem dimension when selecting the vector y. Since we can always select an  $\mathbf{A}$  with  $\sigma_1(\mathbf{A}) = \sqrt{L}$  and  $\sigma_n(\mathbf{A}) = \sqrt{\mu}$  and a y to ensure (158), we have  $\mathbf{f} \in \mathcal{F}_n$ . It remains to show how to set the right singular vectors of  $\mathbf{A}$  to ensure (154).

Plugging (157) into (8) yields

$$\boldsymbol{x}_{t+1} = \boldsymbol{x}_t - \eta \mathbf{A}^{\top} \mathbf{A} (\boldsymbol{x}_t - \boldsymbol{x}_f^*). \tag{159}$$

Subtracting  $x_f^*$  from both sides of (159), we conclude that the distance to the optimizer  $x_f^*$  satisfies

$$\|\boldsymbol{x}_{t+1} - \boldsymbol{x}_{f}^{*}\| \leq \sigma_{1} \left( \mathbf{I} - \eta \mathbf{A}^{\top} \mathbf{A} \right) \|\boldsymbol{x}_{t} - \boldsymbol{x}_{f}^{*}\|,$$
 (160)

where equality is achieved if  $x_t - x_f^*$  points in the direction corresponding to the largest singular vector of the matrix  $\mathbf{I} - \eta \mathbf{A}^{\top} \mathbf{A}$ . Since

$$\sigma_1 \left( \mathbf{I} - \eta \mathbf{A}^{\top} \mathbf{A} \right) = \max \left\{ \left| 1 - \eta \sigma_n^2(\mathbf{A}) \right|, \left| 1 - \eta \sigma_1^2(\mathbf{A}) \right| \right\},$$
(161)

we designate the unit vector

$$\boldsymbol{v}_1 \triangleq \frac{\boldsymbol{x}_0 - \boldsymbol{x}_{\mathsf{f}}^*}{\|\boldsymbol{x}_0 - \boldsymbol{x}_{\mathsf{f}}^*\|} \tag{162}$$

as the right singular vector of  $\mathbf{A}$  corresponding to either  $\sigma_1(\mathbf{A})$  if  $\left|1-\eta\sigma_1^2(\mathbf{A})\right|$  achieves the maximum in (161) or  $\sigma_n(\mathbf{A})$  otherwise. This is determined solely by the stepsize  $\eta$ ; the optimal stepsize (20) ensures that  $\left|1-\eta\sigma_1^2(\mathbf{A})\right|=\left|1-\eta\sigma_n^2(\mathbf{A})\right|=\sigma_{\mathrm{GD}}$ . To complete the construction of  $\mathbf{A}$ , we complement  $v_1$  with n-1 orthonormal vectors to form an orthonormal basis  $\{v_i\}_{i=1}^n$  of  $\mathbb{R}^n$ . Then,

$$\mathbf{x}_1 - \mathbf{x}_f^* = \sigma_1 \left( \mathbf{I} - \eta \mathbf{A}^\top \mathbf{A} \right) \left( \mathbf{x}_0 - \mathbf{x}_f^* \right),$$
 (163)

and using (159) it is easy to show by induction that

$$\boldsymbol{x}_{t+1} - \boldsymbol{x}_{f}^{*} = \sigma_{1} \left( \mathbf{I} - \eta \mathbf{A}^{\top} \mathbf{A} \right) \left( \boldsymbol{x}_{t} - \boldsymbol{x}_{f}^{*} \right), \tag{164}$$

which implies that (160) holds with equality  $\forall t = 0, 1, \dots$ 

Remark C.1. While variants of Lemma C.1 are known in the literature (e.g. [75, Example 1.3]), they are not directly applicable because of our need to satisfy (15). Furthermore, Lemma C.1 constructs a worst-case problem instance for any initial point and constant stepsize chosen by the GD, whereas the worst-case problem instance constructed in [75, Example 1.3] is tailored to a particular starting point  $x_0 \neq 0$  and the optimal  $\eta$  (20).

## B. Proof of Theorem IV.1

On one hand, we have

$$\inf_{\mathbf{A} \in \mathcal{A}_{\mathrm{GD}}} \sigma_{\mathbf{A}}(n, R) \ge \inf_{\mathbf{A} \in \mathcal{A}_{\mathrm{GD}}} \sigma_{\mathbf{A}}(n, \infty) \tag{165}$$

$$\geq \sigma_{\rm GD},$$
 (166)

where (165) holds because the left side of (165) is non-increasing in the data rate R by definition (1), and (166) is by Lemma C.1 since the infinite-rate algorithm in  $\mathcal{A}_{\rm GD}$  is the one incurring no quantization error at each iteration, i.e., the GD itself.

On the other hand, to show

$$\inf_{\mathbf{A} \in \mathcal{A}_{GD}} \sigma_{\mathbf{A}}(n, R) \ge 2^{-R},\tag{167}$$

we fix an algorithm  $A \in \mathcal{A}_{GD}$  operating at  $R' \leq R$  bits per dimension. The set of possible outputs of A after T iterations

$$S_{A} \triangleq \left\{ \hat{\boldsymbol{x}}_{t} \in \mathbb{R}^{n} \mid \hat{\boldsymbol{x}}_{t} \text{ is the output} \right.$$
of A after T iterations \} (168)

has cardinality

$$|\mathcal{S}_{\mathsf{A}}| = 2^{nR'T}.\tag{169}$$

Given  $S_A$ , consider the minimum-distance quantizer  $q_A$ 

$$q_{A}(\boldsymbol{x}) = \operatorname*{arg\,min}_{\hat{\boldsymbol{x}} \in \mathcal{S}_{A}} \|\hat{\boldsymbol{x}} - \boldsymbol{x}\| \tag{170}$$

with dynamic range D and covering radius  $d(q_A)$  (18). In other words,  $2^{nR'T}$  Euclidean balls of radius  $d(q_A)$  with centers in  $S_A$  cover  $\mathcal{B}(D)$ ; therefore

$$\frac{D}{\mathsf{d}\left(\mathsf{q}_{\mathsf{A}}\right)} \le 2^{R'T}.\tag{171}$$

(This classical volume-division argument also shows that  $\rho(q_A) \ge 1$  (19)).

Since one can construct an  $f \in \mathcal{F}_n$  such that

$$\|\mathbf{q}_{A}(x_{f}^{*}) - x_{f}^{*}\| = \mathsf{d}(\mathbf{q}_{A}),$$
 (172)

(167) follows by rearranging (171) and taking the limit  $T \rightarrow \infty$ .

#### C. Proof of Theorem IV.2

The following lemma shows that gradient iterative methods cannot achieve an arbitrarily low contraction factor.

**Lemma C.2** ([66, Th. 2.1.13]). For any gradient method (54), there exists an L-smooth and  $\mu$ -strongly convex function  $f: \mathbb{L}_2 \to \mathbb{R}$  such that,  $\forall t = 0, 1, \ldots$ ,

$$\|x_t - x_f^*\| \ge \sigma_{HB}^t \|x_0 - x_f^*\|.$$
 (173)

The proof of Theorem IV.2 follows the same steps as the proof of Theorem IV.1, with the replacement of the converse for unquantized GD (Lemma C.1) by Lemma C.2.

#### APPENDIX D

#### CONVERGENCE ANALYSIS OF K-WORKER NQ-GD

The path of NQ-GD satisfies the following recursive relation.

**Lemma D.1** (NQ-GD trajectory). At each iteration t = 0, 1, 2, ..., the path of NQ-GD with stepsize (20) satisfies

$$\|\hat{\boldsymbol{x}}_{t+1} - \boldsymbol{x}_{f}^{*}\| \le \sigma_{\text{GD}} \|\hat{\boldsymbol{x}}_{t} - \boldsymbol{x}_{f}^{*}\| + \frac{\eta_{\text{GD}}}{K} \sum_{k=1}^{K} \|\boldsymbol{e}_{t,k}\|, \quad (174)$$

where  $e_{t,k} \triangleq q_{t,k} - \nabla f_k(\hat{x}_t)$ , and  $\sigma_{GD}$  is defined in (21).

*Proof of Lemma D.1.* The update rule (Line 5) and the quantizer inputs (Line 4) together imply

$$\hat{x}_{t+1} = \hat{x}_t - \frac{\eta_{\text{GD}}}{K} \sum_{k=1}^{K} (\nabla f_k(\hat{x}_t) + e_{t,k})$$
 (175)

$$= \hat{\boldsymbol{x}}_t - \eta_{\text{GD}} \nabla f(\hat{\boldsymbol{x}}_t) - \frac{\eta_{\text{GD}}}{K} \sum_{k=1}^K \boldsymbol{e}_{t,k}$$
 (176)

where (176) is due to (59). Triangle inequality now implies

$$\|\hat{\boldsymbol{x}}_{t+1} - \boldsymbol{x}_{\mathsf{f}}^*\| \tag{177}$$

$$\leq \|\hat{x}_{t} - \eta_{\text{GD}} \nabla f(\hat{x}_{t}) - x_{f}^{*}\| + \frac{\eta_{\text{GD}}}{K} \sum_{k=1}^{K} \|e_{t,k}\|.$$
 (178)

Applying the coercive property of smooth and strongly convex functions [66, Th. 2.1.12] in the same manner as in [66, Proof of Th. 2.1.5] to further upper-bound the first term in (178) gives (174).

Denote for brevity

$$\sigma \triangleq \sigma_{\rm GD} + C,\tag{179}$$

$$C \triangleq \frac{\eta_{\text{GD}}\rho_n}{K} \sum_{k=1}^K \frac{L_k}{2^{R_k}}.$$
 (180)

Minimizing (179) under the sum rate constrant (60) is a convex optimization problem whose solution is given by (67) and whose optimal value is given by the right side of (69). The following nonasymptotic result, which with the optimal rate allocation immediately yields Theorem V.1, uses an inductive argument to simultaneously bound both terms in the right-hand side of (174).

**Theorem D.1** (Convergence of NQ-GD). In the setting of Theorem V.1,

$$\|\hat{\boldsymbol{x}}_t - \boldsymbol{x}_f^*\| \le \sigma^t D,\tag{181}$$

where  $\sigma$  is defined in (179).

*Proof.* We prove (181) by induction.

- Base case: for t = 0, (181) holds by (15).
- Inductive step: Suppose (181) holds for iteration t. Then, by  $L_k$ -smoothness of  $f_k$ , (62) and the inductive hypothesis,

$$\|\nabla f_k(\hat{\boldsymbol{x}}_t)\| \le L_k \|\hat{\boldsymbol{x}}_t - \boldsymbol{x}_f^*\| \tag{182}$$

$$\leq \sigma^t L_k D,\tag{183}$$

which due to the choice of dynamic ranges (66) and (26) leads to

$$\|\boldsymbol{e}_{t,k}\| \le \frac{\rho_n}{2^{R_k}} \sigma^t L_k D. \tag{184}$$

Applying (184) and the inductive hypothesis to (174) yields

$$\|\hat{x}_{t+1} - x_{f}^{*}\| < \sigma^{t} [\sigma_{GD} + C] D,$$
 (185)

which is exactly (181) for t + 1 if the optimal rate allocation is employed.

Without assumption (iv) (62) that the summands  $f_k$  share the minimizer, NQ-GD converges only to a neighborhood of  $x_f^*$ , albeit exponentially fast. The radius of this neighborhood vanishes as the data rate  $R \to \infty$ :

**Theorem D.2.** In the setting of Theorem V.1 but without assumption (iv) (62) on the objective function, setting the dynamic ranges to  $r'_{t,k}$ , where

$$r'_{t,k} = r_{t,k} + 2L_k D\left(C\sum_{\tau=0}^{t-1} \sigma^{\tau} + 1\right),$$
 (186)

 $r_{t,k}$  is defined in (66), and C is defined in (180), NQ-GD achieves

$$\|\hat{\boldsymbol{x}}_t - \boldsymbol{x}_f^*\| \le \sigma^t D + 2DC \sum_{\tau=0}^{t-1} \sigma^{\tau},$$
(187)

where  $\sigma$  is defined in (179).

*Proof.* We follow the reasoning in the proof of Theorem D.1 until (182), which we replace by

$$\|\nabla f_k(\hat{\boldsymbol{x}}_t)\| \le L_k \|\hat{\boldsymbol{x}}_t - \boldsymbol{x}_{f_k}^*\| \tag{188}$$

$$\leq L_k \left( \left\| \hat{\boldsymbol{x}}_t - \boldsymbol{x}_f^* \right\| + \left\| \boldsymbol{x}_{f_k}^* - \boldsymbol{x}_f^* \right\| \right) \tag{189}$$

$$\leq L_k (\|\hat{x}_t - x_f^*\| + 2D)$$
 (190)

$$\leq L_k \left( \sigma^t D + 2DC \sum_{\tau=0}^{t-1} \sigma^{\tau} + 2D \right), \quad (191)$$

where (191) applies inductive hypothesis (187). Due to (191), the choice of dynamic ranges (186) ensures that the quantizer input stays within its dynamic range, which limits the quantization error to  $\rho_n 2^{-R_k}$  times the right-hand side of (191) (recall (26)). Plugging this bound on quantization error in (174) and applying inductive hypothesis (187) again establishes (187) for t+1.

#### APPENDIX E

ANALYSIS OF QUANTIZED GD WITH ERROR FEEDBACK OF [10], [22], [44]

We derive the following convergence guarantee on quantized GD with the error feedback mechanism of [10], [22], [44].

**Lemma E.1** (Trajectory of GD with error feedback of [10], [22], [44]). Let f be an L-smooth and  $\mu$ -strongly convex function on  $\mathbb{R}^n$ . Then, the distance to the optimizer at each iteration  $t \in \mathbb{N}$  of quantized gradient descent (50) with quantizer input (30) and stepsize (20) is bounded as

$$\|\hat{x}_{t+1} - x^*\| \le \sigma_{GD} \|\hat{x}_t - x^*\| + \eta_{GD} \|e_t\| + (\sigma_{GD} + \eta_{GD} L) \eta \|e_{t-1}\|.$$
 (192)

*Proof.* Using (30), the quantizer output can be computed as

$$\mathbf{q}_t = \left[\nabla f(\hat{\mathbf{x}}_t) + \mathbf{e}_t\right] - \mathbf{e}_{t-1},\tag{193}$$

and plugging (193) into (50) gives

$$\hat{\boldsymbol{x}}_{t+1} = \hat{\boldsymbol{x}}_t - \eta \nabla f(\hat{\boldsymbol{x}}_t) - \eta \boldsymbol{e}_t + \eta \boldsymbol{e}_{t-1}. \tag{194}$$

Denoting a shifted trajectory by

$$\tilde{\boldsymbol{x}}_t \triangleq \hat{\boldsymbol{x}}_t + \eta \boldsymbol{e}_{t-1},\tag{195}$$

we rewrite (194) as

$$\tilde{\boldsymbol{x}}_{t+1} = \tilde{\boldsymbol{x}}_t - \eta \nabla f(\tilde{\boldsymbol{x}}_t - \eta \boldsymbol{e}_{t-1}), \tag{196}$$

which via triangle inequality implies

$$\|\tilde{\boldsymbol{x}}_{t+1} - \boldsymbol{x}^*\| \le \|\tilde{\boldsymbol{x}}_t - \boldsymbol{x}^* - \eta \nabla f(\tilde{\boldsymbol{x}}_t)\|$$

$$\tag{197}$$

$$+ \eta \left\| \nabla \mathsf{f}(\tilde{\boldsymbol{x}}_t) - \nabla \mathsf{f}(\tilde{\boldsymbol{x}}_t - \eta \boldsymbol{e}_{t-1}) \right\| \quad (198)$$

$$<\sigma_{GD} \|\tilde{\boldsymbol{x}}_{t} - \boldsymbol{x}^*\| + \eta^2 L \|\boldsymbol{e}_{t-1}\|,$$
 (199)

where the first term in the right side of (199) is obtained by the convergence guarantee of GD (Lemma B.2), and the second term is due to L-smoothness (13). Since (195) implies

$$\|\tilde{x}_t - x^*\| - \eta \|e_{t-1}\| \le \|\hat{x}_t - x^*\|$$
 (200)

$$\leq \|\tilde{\boldsymbol{x}}_t - \boldsymbol{x}^*\| + \eta \|\boldsymbol{e}_{t-1}\|, \quad (201)$$

we leverage (199) to control  $\|\hat{x}_{t+1} - x^*\|$  as follows:

$$\|\hat{x}_{t+1} - x^*\|$$

$$\leq \sigma_{GD} \|\tilde{x}_t - x^*\| + \eta^2 L \|e_{t-1}\| + \eta \|e_t\|$$

$$\leq \sigma_{GD} \|\hat{x}_t - x^*\| + \eta \|e_t\| + (\sigma_{GD} + \eta L)\eta \|e_{t-1}\|,$$
 (202)

where (203) is due to (200).

Compared to the guarantee on NQ-GD in Lemma D.1, that on GD with error feedback of [10], [22], [44] in Lemma E.1 has an extra error term. Thus, it is unclear whether the error feedback mechanism of [10], [22], [44] can even improve upon NQ-GD in our nonstochastic problem setting.

#### REFERENCES

- C.-Y. Lin, V. Kostina, and B. Hassibi, "Differentially quantized gradient descent," in *Proceedings 2021 IEEE International Symposium on Information Theory*, July 2021.
- [2] M. Zinkevich, M. Weimer, L. Li, and A. J. Smola, "Parallelized stochastic gradient descent," in *Advances in Neural Information Processing Systems* 23, Vancouver, British Columbia, Canada, Dec. 2010, pp. 2595–2603
- [3] B. Recht, C. Re, S. Wright, and F. Niu, "Hogwild: A lock-free approach to parallelizing stochastic gradient descent," in *Advances in Neural Information Processing Systems* 24, Granada, Spain, Dec. 2011, pp. 693–701.
- [4] R. Bekkerman, M. Bilenko, and J. Langford, Scaling Up Machine Learning: Parallel and Distributed Approaches. New York, NY, USA: Cambridge University Press, Dec. 2011.
- [5] J. Dean, G. Corrado, R. Monga, K. Chen, M. Devin, M. Mao, M. Ranzato, A. Senior, P. Tucker, K. Yang, Q. V. Le, and A. Y. Ng, "Large scale distributed deep networks," in *Advances in Neural Information Processing Systems* 25, Lake Tahoe, NV, USA, Dec. 2012, pp. 1223–1231.
- [6] T. Chilimbi, Y. Suzue, J. Apacible, and K. Kalyanaraman, "Project Adam: Building an efficient and scalable deep learning training system," in 11th USENIX Symposium on Operating Systems Design and Implementation (OSDI 14), Broomfield, CO, Oct. 2014, pp. 571–582.
- [7] C. M. De Sa, C. Zhang, K. Olukotun, C. Ré, and C. Ré, "Taming the wild: A unified analysis of hogwild-style algorithms," in *Advances in Neural Information Processing Systems* 28, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, Eds. Curran Associates, Dec. 2015, pp. 2674–2682.
- [8] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon, "Federated learning: Strategies for improving communication efficiency," in NIPS Workshop on Private Multi-Party Machine Learning, Dec. 2016.
- [9] K. Scaman, F. Bach, S. Bubeck, Y. T. Lee, and L. Massoulié, "Optimal algorithms for smooth and strongly convex distributed optimization in networks," in *Proceedings of the 34th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, D. Precup and Y. W. Teh, Eds., vol. 70. International Convention Centre, Sydney, Australia: PMLR, 06–11 Aug 2017, pp. 3027–3036.
- [10] F. Seide, H. Fu, J. Droppo, G. Li, and D. Yu, "1-bit stochastic gradient descent and application to data-parallel distributed training of speech DNNs," in *Interspeech 2014*, Sep. 2014.
- [11] M. Li, D. G. Andersen, A. J. Smola, and K. Yu, "Communication efficient distributed machine learning with the parameter server," in Advances in Neural Information Processing Systems 27, Montreal, QB, Canada, Dec. 2014, pp. 19–27.
- [12] N. Strom, "Scalable distributed DNN training using commodity GPU cloud computing," in *INTERSPEECH*, Sep. 2015.
- [13] S. Zhang, A. E. Choromanska, and Y. LeCun, "Deep learning with elastic averaging SGD," in *Advances in Neural Information Processing Systems* 28, Montreal, QB, Canada, Dec. 2015, pp. 685–693.
- [14] H. Robbins and S. Monro, "A stochastic approximation method," Ann. Math. Statist., vol. 22, no. 3, pp. 400–407, Sep. 1951.
- [15] 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," in *Advances in Neural Information Processing Systems 30*, Long Beach, CA, USA, Dec. 2017, pp. 1509–1519.

- [16] J. Bernstein, Y.-X. Wang, K. Azizzadenesheli, and A. Anandkumar, "signSGD: Compressed optimisation for non-convex problems," in *Proceedings of the 35th International Conference on Machine Learning*, Long Beach, CA, USA, Jul. 2018, pp. 560–569.
- [17] D. Alistarh, D. Grubic, J. Li, R. Tomioka, and M. Vojnovic, "QSGD: Communication-efficient SGD via gradient quantization and encoding," in *Advances in Neural Information Processing Systems 30*, Long Beach, CA, USA, Dec. 2017, pp. 1709–1720.
- [18] A. Ramezani-Kebrya, F. Faghri, and D. M. Roy, "NUQSGD: Improved communication efficiency for data-parallel SGD via nonuniform quantization," Aug. 2019.
- [19] P. Mayekar and H. Tyagi, "Limits on gradient compression for stochastic optimization," arXiv, vol. 2001.09032, Jan. 2020.
- [20] V. Gandikota, D. Kane, R. K. Maity, and A. Mazumdar, "vqSGD: Vector quantized stochastic gradient descent," Nov. 2019.
- [21] A. F. Aji and K. Heafield, "Sparse communication for distributed gradient descent," in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, Copenhagen, Denmark, Sep. 2017, pp. 440–445.
- [22] S. U. Stich, J.-B. Cordonnier, and M. Jaggi, "Sparsified SGD with memory," in *Advances in Neural Information Processing Systems 31*, Montréal, Canada, Dec. 2018, pp. 4447–4458.
- [23] J. Wangni, J. Wang, J. Liu, and T. Zhang, "Gradient sparsification for communication-efficient distributed optimization," in *Advances in Neural Information Processing Systems 31*, Montréal, Canada, Dec. 2018, pp. 1299–1309.
- [24] H. Wang, S. Sievert, S. Liu, Z. Charles, D. Papailiopoulos, and S. Wright, "ATOMO: Communication-efficient learning via atomic sparsification," in *Advances in Neural Information Processing Systems* 31, Montréal, Canada, Dec. 2018, pp. 9850–9861.
- [25] Y. Lin, S. Han, H. Mao, Y. Wang, and B. Dally, "Deep gradient compression: Reducing the communication bandwidth for distributed training," in *International Conference on Learning Representations*, Vancouver, BC, Canada, Apr. 2018.
- [26] N. Dryden, S. A. Jacobs, T. Moon, and B. Van Essen, "Communication quantization for data-parallel training of deep neural networks," in Proceedings of the Workshop on Machine Learning in High Performance Computing Environments, 2016, pp. 1–8.
- [27] D. Alistarh, T. Hoefler, M. Johansson, N. Konstantinov, S. Khirirat, and C. Renggli, "The convergence of sparsified gradient methods," in *Advances in Neural Information Processing Systems 31*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds. Curran Associates, Dec. 2018, pp. 5973–5983.
- [28] M. Yu, Z. Lin, K. Narra, S. Li, Y. Li, N. S. Kim, A. Schwing, M. Annavaram, and S. Avestimehr, "GradiVeQ: Vector quantization for bandwidth-efficient gradient aggregation in distributed CNN training," *Advances in Neural Information Processing Systems*, vol. 31, pp. 5123– 5133, Jan. 2018.
- [29] M. M. Amiri and D. Gündüz, "Machine learning at the wireless edge: Distributed stochastic gradient descent over-the-air," in 2019 IEEE International Symposium on Information Theory (ISIT), 2019, pp. 1432– 1436.
- [30] G. Zhu, Y. Wang, and K. Huang, "Broadband analog aggregation for low-latency federated edge learning," *IEEE Transactions on Wireless Communications*, vol. 19, no. 1, pp. 491–506, 2020.
- [31] K. Yang, T. Jiang, Y. Shi, and Z. Ding, "Federated learning via overthe-air computation," *IEEE Transactions on Wireless Communications*, vol. 19, no. 3, pp. 2022–2035, 2020.
- [32] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro, "Robust stochastic approximation approach to stochastic programming," SIAM Journal on Optimization, vol. 19, no. 4, pp. 1574–1609, Jan. 2009.
- [33] S. Bubeck, "Convex optimization: Algorithms and complexity," Found. Trends Mach. Learn., vol. 8, no. 3-4, pp. 231–357, Nov. 2015.
- [34] L. Bottou, F. E. Curtis, and J. Nocedal, "Optimization methods for large-scale machine learning," SIAM Review, vol. 60, no. 2, pp. 223–311, May 2018.
- [35] K. Mishchenko, E. Gorbunov, M. Takáč, and P. Richtárik, "Distributed learning with compressed gradient differences," Jan. 2019.
- [36] S. Horváth, D. Kovalev, K. Mishchenko, S. Stich, and P. Richtárik, "Stochastic distributed learning with gradient quantization and variance reduction," Apr. 2019.
- [37] S. Horváth, C.-Y. Ho, L. Horváth, A. N. Sahu, M. Canini, and P. Richtárik, "Natural compression for distributed deep learning," May. 2019
- [38] C. Philippenko and A. Dieuleveut, "Bidirectional compression in heterogeneous settings for distributed or federated learning with partial partici-

- pation: tight convergence guarantees," arXiv preprint arXiv:2006.14591, June 2020.
- [39] E. Gorbunov, F. Hanzely, and P. Richtárik, "A unified theory of SGD: Variance reduction, sampling, quantization and coordinate descent," in International Conference on Artificial Intelligence and Statistics, June 2020, pp. 680–690.
- [40] E. Gorbunov, D. Kovalev, D. Makarenko, and P. Richtárik, "Linearly converging error compensated SGD," *Advances in Neural Information Processing Systems*, vol. 33, pp. 20889–20900, Dec. 2020.
- [41] A. Khaled, O. Sebbouh, N. Loizou, R. M. Gower, and P. Richtárik, "Unified analysis of stochastic gradient methods for composite convex and smooth optimization," arXiv preprint arXiv:2006.11573, June 2020.
- [42] E. Gorbunov, K. P. Burlachenko, Z. Li, and P. Richtárik, "MARINA: Faster non-convex distributed learning with compression," in *International Conference on Machine Learning*. PMLR, July 2021, pp. 3788–3798.
- [43] R. Islamov, X. Qian, and P. Richtárik, "Distributed second order methods with fast rates and compressed communication," in *International Conference on Machine Learning*. PMLR, July 2021, pp. 4617–4628.
- [44] S. P. Karimireddy, Q. Rebjock, S. Stich, and M. Jaggi, "Error feedback fixes SignSGD and other gradient compression schemes," in *Proceedings* of the 36th International Conference on Machine Learning, vol. 97, Long Beach, CA, USA, June 2019, pp. 3252–3261.
- [45] A. Beznosikov, S. Horváth, P. Richtárik, and M. Safaryan, "On biased compression for distributed learning," arXiv:2002.12410, Feb. 2020.
- [46] S. Magnússon, H. Shokri-Ghadikolaei, and N. Li, "On maintaining linear convergence of distributed learning and optimization under limited communication," *IEEE Transactions on Signal Processing*, vol. 68, pp. 6101–6116, 2020.
- [47] R. M. Gray, Source Coding Theory. Kluwer Academic Publishers, Oct. 1989
- [48] S. Zheng, Z. Huang, and J. Kwok, "Communication-efficient distributed blockwise momentum SGD with error-feedback," in *Advances in Neural Information Processing Systems 32*. Curran Associates, Dec. 2019, pp. 11 450–11 460.
- [49] J. Wu, W. Huang, J. Huang, and T. Zhang, "Error compensated quantized SGD and its applications to large-scale distributed optimization," in Proceedings of the 35th International Conference on Machine Learning, vol. 80, Stockholm, Sweden, July 2018, pp. 5325–5333.
- [50] J. Sun, T. Chen, G. B. Giannakis, Q. Yang, and Z. Yang, "Lazily aggregated quantized gradient innovation for communication-efficient federated learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [51] X. Qian, P. Richtárik, and T. Zhang, "Error compensated distributed sgd can be accelerated," Advances in Neural Information Processing Systems, vol. 34, Dec. 2021.
- [52] P. Richtárik, I. Sokolov, and I. Fatkhullin, "EF21: A new, simpler, theoretically better, and practically faster error feedback," Advances in Neural Information Processing Systems, vol. 34, Dec. 2021.
- [53] S. Horváth and P. Richtárik, "A better alternative to error feedback for communication-efficient distributed learning," in *Eighth International Conference on Learning Representations*, Apr. 2020.
- [54] J. Acharya, C. De Sa, D. Foster, and K. Sridharan, "Distributed learning with sublinear communication," in *Proceedings of the 36th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, K. Chaudhuri and R. Salakhutdinov, Eds., vol. 97. Long Beach, CA, USA: PMLR, 09–15 Jun. 2019, pp. 40–50.
- [55] P. Mayekar and H. Tyagi, "RATQ: A universal fixed-length quantizer for stochastic optimization," arXiv, vol. 1908.08200, Dec. 2019.
- [56] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Found. Trends Mach. Learn.*, vol. 3, no. 1, pp. 1–122, Jan 2011.
- [57] Y. Zhang, J. C. Duchi, and M. J. Wainwright, "Communication-efficient algorithms for statistical optimization," *J. Mach. Learn. Res.*, vol. 14, no. 1, pp. 3321–3363, Jan 2013.
- [58] A. T. Suresh, F. X. Yu, S. Kumar, and H. B. McMahan, "Distributed mean estimation with limited communication," in *Proceedings of the* 34th International Conference on Machine Learning, ser. Proceedings of Machine Learning Research, vol. 70. PMLR, 06–11 Aug 2017, pp. 3329–3337.
- [59] A. Nedic, A. Olshevsky, A. Ozdaglar, and J. N. Tsitsiklis, "On distributed averaging algorithms and quantization effects," *IEEE Transactions on Automatic Control*, vol. 54, no. 11, pp. 2506–2517, Nov 2009.
- [60] A. Reisizadeh, A. Mokhtari, H. Hassani, and R. Pedarsani, "An exact quantized decentralized gradient descent algorithm," *IEEE Transactions* on Signal Processing, vol. 67, no. 19, pp. 4934–4947, Oct 2019.

- [61] Y. Zhang, J. Duchi, M. I. Jordan, and M. J. Wainwright, "Information-theoretic lower bounds for distributed statistical estimation with communication constraints," in *Advances in Neural Information Processing Systems* 26, 2013, pp. 2328–2336.
- [62] J. N. Tsitsiklis and Z.-Q. Luo, "Communication complexity of convex optimization," *Journal of Complexity*, vol. 3, no. 3, pp. 231 – 243, 1987.
- [63] Y. Arjevani and O. Shamir, "Communication complexity of distributed convex learning and optimization," in *Advances in Neural Information Processing Systems* 28, Dec. 2015, pp. 1756–1764.
- [64] M. Li, D. G. Andersen, J. W. Park, A. J. Smola, A. Ahmed, V. Josifovski, J. Long, E. J. Shekita, and B.-Y. Su, "Scaling distributed machine learning with the parameter server," in 11th USENIX Symposium on Operating Systems Design and Implementation, Broomfield, CO, Oct. 2014, pp. 583–598.
- [65] S. Khirirat, H. R. Feyzmahdavian, and M. Johansson, "Distributed learning with compressed gradients," arXiv: 1806.06573, June 2018.
- [66] Y. Nesterov, Introductory Lectures on Convex Optimization: A Basic Course. Springer, Dec. 2014.
- [67] B. T. Polyak, Introduction to optimization, ser. Translations Series in Mathematics and Engineering. Optimization Software, May 1987.
- [68] C. A. Rogers, "Covering a sphere with spheres," *Mathematika*, vol. 10, no. 2, pp. 157—164, Dec. 1963.
- [69] M. Friedlander and M. Schmidt, "Hybrid deterministic-stochastic methods for data fitting," SIAM Journal on Scientific Computing, vol. 34, no. 3, pp. A1380–A1405, May 2012.
- [70] Y. Nesterov, "A method of solving a convex programming problem with convergence rate  $O(1/k^2)$ ," in *Soviet Mathematics Doklady*, vol. 27, no. 2, 1982, pp. 372–376.
- [71] Z. Allen-Zhu and L. Orecchia, "Linear coupling: An ultimate unification of gradient and mirror descent," arXiv: 1407.1537, July 2014.
- [72] W. Su, S. Boyd, and E. J. Candès, "A differential equation for modeling Nesterov's accelerated gradient method: Theory and insights," *Journal of Machine Learning Research*, vol. 17, no. 153, pp. 1–43, 2016.
- [73] R. Zamir, Lattice Coding for Signals and Networks: A Structured Coding Approach to Quantization, Modulation, and Multiuser Information Theory. Cambridge University Press, Sep. 2014.
- [74] S. Kolodziej, M. Aznaveh, M. Bullock, J. David, T. Davis, M. Henderson, Y. Hu, and R. Sandstrom, "The SuiteSparse matrix collection website interface," *Journal of Open Source Software*, vol. 4, no. 35, p. 1244, Mar. 2019.
- [75] E. de Klerk, F. Glineur, and A. B. Taylor, "On the worst-case complexity of the gradient method with exact line search for smooth strongly convex functions," *Optim. Lett.*, vol. 11, no. 7, pp. 1185–1199, Oct. 2017.
- [76] Y. Arjevani, S. Shalev-Shwartz, and O. Shamir, "On lower and upper bounds in smooth and strongly convex optimization," *Journal of Ma*chine Learning Research, vol. 17, no. 126, pp. 1–51, 2016.
- [77] S. Boyd and L. Vandenberghe, Convex Optimization. Cambridge University Press, 2004.
- [78] Y. S. Abu-Mostafa, M. Magdon-Ismail, and H.-T. Lin, *Learning From Data*. AMLBook, 2012.
- [79] M. Schmidt and N. L. Roux, "Fast convergence of stochastic gradient descent under a strong growth condition." Aug. 2013.
- descent under a strong growth condition," Aug. 2013.
  [80] D. Needell, R. Ward, and N. Srebro, "Stochastic gradient descent, weighted sampling, and the randomized Kaczmarz algorithm," in *Advances in Neural Information Processing Systems*, vol. 27. Curran Associates, Dec. 2014, pp. 1017–1025.
- [81] S. Ma, R. Bassily, and M. Belkin, "The power of interpolation: Understanding the effectiveness of SGD in modern over-parametrized learning," in *International Conference on Machine Learning*, vol. 80, Stockholm, Sweden, July 2018, pp. 3325–3334.
- [82] F. Wirth, "On the calculation of time-varying stability radii," *International Journal of Robust and Nonlinear Control: IFAC-Affiliated Journal*, vol. 8, no. 12, pp. 1043–1058, 1998.

Chung-Yi Lin received the B.Sc. degree in electrical engineering and the M.Sc. degree in communication engineering from the National Taiwan University in 2015 and 2018, respectively. He was a Ph.D. student in electrical engineering at the California Institute of Technology between 2018 and 2020. He has been a quantitative researcher at the Kronos Research based in Taiwan since 2021.

**Victoria Kostina** (S'12–M'14) received the bachelor's degree from Moscow Institute of Physics and Technology (MIPT) in 2004, the master's degree from University of Ottawa in 2006, and the Ph.D. degree from Princeton University in 2013. During her studies at MIPT, she was affiliated with the Institute for Information Transmission Problems of the Russian Academy of Sciences.

She is currently a Professor of electrical engineering and computing and mathematical sciences at California Institute of Technology. Her research interests include information theory, coding, control, learning, and communications. She received the Natural Sciences and Engineering Research Council of Canada postgraduate scholarship during 2009–2012, Princeton Electrical Engineering Best Dissertation Award in 2013, Simons-Berkeley research fellowship in 2015 and the NSF CAREER award in 2017.

She is a Guest Editor for the IEEE Journal on Selected Areas in Information Theory 2022 special issue on Modern Compression.

**Babak Hassibi** (Member, IEEE) was born in Tehran, Iran, in 1967. He received the B.S. degree from the University of Tehran in 1989, and the M.S. and Ph.D. degrees from Stanford University in 1993 and 1996, respectively, all in electrical engineering.

He has been with the California Institute of Technology since January 2001, where he is currently the Mose and Lilian S. Bohn Professor of Electrical Engineering. From 2013-2016 he was the Gordon M. Binder/Amgen Professor of Electrical Engineering and from 2008-2015 he was Executive Officer of Electrical Engineering, as well as Associate Director of Information Science and Technology. From October 1996 to October 1998 he was a research associate at the Information Systems Laboratory, Stanford University, and from November 1998 to December 2000 he was a Member of the Technical Staff in the Mathematical Sciences Research Center at Bell Laboratories. Murray Hill, NJ. He has also held short-term appointments at Ricoh California Research Center, the Indian Institute of Science, and Linkoping University, Sweden. His research interests include communications and information theory, control and network science, and signal processing and machine learning. He is the coauthor of the books (both with A.H. Sayed and T. Kailath) Indefinite Quadratic Estimation and Control: A Unified Approach to H<sup>2</sup> and  $H^{\infty}$  Theories (New York: SIAM, 1999) and Linear Estimation (Englewood Cliffs, NJ: Prentice Hall, 2000). He is a recipient of an Alborz Foundation Fellowship, the 1999 O. Hugo Schuck best paper award of the American Automatic Control Council (with H. Hindi and S.P. Boyd), the 2002 National Science Foundation Career Award, the 2002 Okawa Foundation Research Grant for Information and Telecommunications, the 2003 David and Lucille Packard Fellowship for Science and Engineering, the 2003 Presidential Early Career Award for Scientists and Engineers (PECASE), and the 2009 Al-Marai Award for Innovative Research in Communications, and was a participant in the 2004 National Academy of Engineering "Frontiers in Engineering" program.

He has been a Guest Editor for the IEEE Transactions on Information Theory special issue on "space-time transmission, reception, coding and signal processing" was an Associate Editor for Communications of the IEEE Transactions on Information Theory during 2004-2006, and is currently an Editor for the Journal "Foundations and Trends in Information and Communication" and for the IEEE Transactions on Network Science and Engineering. He is an IEEE Information Theory Society Distinguished Lecturer for 2016-2017 and was General Co-Chair if the 2020 IEEE International Symposium on Information Theory (ISIT 2020).