
DoStoVoQ: Doubly Stochastic Voronoi Vector Quantization SGD for Federated Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The growing size of models and datasets have made distributed implementation
2 of stochastic gradient descent (SGD) an active field of research. However the
3 high bandwidth cost of communicating gradient updates between nodes remains
4 a bottleneck; lossy compression is a way to alleviate this problem. We propose a
5 new *unbiased* Vector Quantizer (VQ), named StoVoQ, to perform gradient quanti-
6 zation. This approach relies on introducing randomness within the quantization
7 process, that is based on the use of unitarily invariant random codebooks and on
8 a straightforward bias compensation method. The distortion of StoVoQ signif-
9 icantly improves upon existing quantization algorithms. Next, we explain how
10 to combine this quantization scheme within a Federated Learning framework for
11 complex high-dimensional model (dimension $> 10^6$), introducing DoStoVoQ. We
12 provide theoretical guarantees on the quadratic error and (absence of) bias of the
13 compressor, that allow to leverage strong theoretical results of convergence, e.g.,
14 with heterogeneous workers or variance reduction. Finally, we show that training
15 on convex and non-convex deep learning problems, our method leads to significant
16 reduction of bandwidth use while preserving model accuracy.

17 1 Introduction

18 In this paper, we consider the Federated Learning framework, in which a potentially large number K
19 of *workers* cooperate to solve the following problem:

$$\min_{\theta \in \mathbb{R}^D} \sum_{k=1}^K f_k(\theta), \quad (1)$$

20 where each function $f_k : \mathbb{R}^D \rightarrow \mathbb{R}$ represents the empirical risk on worker $k \in [K]$ (where
21 $[K] = \{1, \dots, K\}$) and D is the ambient dimension of our problem. Each worker potentially holds a
22 fraction of the data, and can share information with a central server, which progressively aggregates
23 and updates the model accordingly [18, 17].

24 Stochastic gradient algorithms [28] are particularly well suited in the *large scale learning* setting [6,
25 7]. The methods can easily be adapted to the distributed (and more generally federated) learning
26 framework; see [17] and the references therein. For synchronous distributed Stochastic Gradient
27 Descent, at every iteration, given the current parameter θ_t , each worker computes an unbiased estimate
28 $g_{k,t+1}(\theta_t)$ of the gradient of the local loss function f_k . The central server then aggregates those
29 oracles and performs the update.

30 Communicating the gradients from the local workers to the central server is often a major bottleneck.
31 The drastic increase both in the number of parameters and of workers over the last years, has made
32 this problem even more acute. Alleviating the communication cost is one of the crucial challenges of

33 federated learning [17, Sec. 3.5]. A central idea to tackle this issue is *communication compression*,
 34 which consists in applying a lossy compression to the parameters or gradients to be transmitted.
 35 Since compression alters the message transmitted, the number of iterations required to reach a given
 36 accuracy may increase, therefore compression is of interest in situations where the communication
 37 gains are large relative to the increase of communication rounds. The design of new compression
 38 schemes (see among others [30, 2, 4, 5, 34]) and the adaptation of the learning algorithms to this
 39 setting (see e.g. [32, 1, 35, 33, 36, 22, 26, 12, 11, 21] and the references therein) are an extremely
 40 active field of research.

41 Our main contribution is to introduce a novel **unbiased vector quantization** procedure allowing to
 42 reach **high-compression rate**, with a **small computational** overhead. More precisely, our contribu-
 43 tions are as follow: first, we introduce StVoQ, a vector quantization algorithm based on unitarily
 44 invariant random codebooks to automatically obtain **directionally unbiased** gradient oracles, and
 45 introduce a scalar **correction function**, that makes compression operator **unbiased** for a very modest
 46 computational cost. We further provide theoretical guarantees on the distortion of the compressor. In
 47 summary, StVoQ algorithm is based on the following points, that are developed in Section 2.

- 48 1. **Vector quantization** The input vector $x \in \mathbb{R}^d$ is mapped onto its nearest neighbor in a codebook
 49 $\mathcal{C}_M = \{c_i\}_{i=1}^M$.
- 50 2. **Random codebook.** A **new codebook** is sampled every time a new quantization operation is
 51 performed. The proposed approach is different from classical random VQ which typically uses a
 52 random codebook, but which is sampled once and then kept fixed.
- 53 3. **Bias removal.** By relying on unitarily invariant distribution for the codewords generation, the
 54 quantized value of each vector $x \in \mathbb{R}^d$ is **directionally unbiased**. The bias only depends on the
 55 number and distributions of the random of codewords and on $\|x\|$. This key property allows to
 56 derive a simple way to remove the quantization bias.

57 Then, we describe how to use StVoQ within the FL framework: this yields the algorithm DoStVoQ.
 58 We prove that this process satisfies a strong assumption on the compression process, that allows to
 59 automatically derive fast convergence rates. In Section 3, we describe DoStVoQ, i.e., how we solve
 60 the optimization problem (1) in dimension D .

- 61 4. **Splitting and renormalizing gradients.** First, we split each gradient to compress into *buckets*
 62 $(x_i)_{i=1, \dots, L}$ of dimension \mathbb{R}^d , to use StVoQ for each bucket.
- 63 5. **Synchronisation of random sequences of codebooks.** We ensure that those codebooks are
 64 independent, at each step and between each machine, by generating a new codebook each time.
 65 To avoid any subsequent communication cost, we synchronously generate the codebooks on the
 66 central and local servers, by initially sharing random seeds.

67 Remark that point 1 was also used in Dai et al. [8]. Points 2 to 3 and 5 are novel ideas that have not
 68 been leveraged in the FL framework. Finally, we demonstrate the effectiveness of random codebook
 69 quantization for gradient compression by extensive experiments in Section 4 on standard benchmarks
 70 like ImageNet or CIFAR10.

71 2 StVoQ algorithm

72 Several compression operators [34, 27, 10, 4, 8, 36, 37] have been introduced recently as bandwidth
 73 reduction for distributed learning became a major challenge. In this section, we first discuss the
 74 importance of unbiasedness of compression operators in Subsection 2.1. We then present the StVoQ
 75 compression scheme in Subsection 2.2. Finally, we compare StVoQ to competing approaches, both
 76 theoretically and empirically on a small scale example with a high compression rate.

77 2.1 Unbiased gradient estimate to mitigate high compression rates

78 We here discuss an important property to mitigate high compression rates in FL settings. A *compression operator* Comp is a (random) mapping on \mathbb{R}^d . Consider the following assumption:

80 **A1 (Unbiased Compression with relatively bounded variance).** A *compression operator* Comp
 81 is unbiased if for any $x \in \mathbb{R}^d$, $\mathbb{E}[\text{Comp}(x)] = x$. It is said to have a ω -bounded relative variance,
 82 for some $\omega > 0$, if it satisfies, for all $x \in \mathbb{R}^d$, $\mathbb{E}[\|\text{Comp}(x) - x\|^2] \leq \omega \|x\|^2$.

83 The most classical compressors, especially Q-SGD and Rand- H satisfy A 1 with different ω , see
 84 Subsection 2.3 and Table 1. On the other hand, some compression operators are biased, i.e.,
 85 $\mathbb{E}[\text{Comp}(x)] \neq x$ for some $x \in \mathbb{R}$. Those operators are often deterministic, as is the case for
 86 Top- H compressor. The most classical assumption for biased operators, is the following contractive
 87 property along the direction of descent [32, 5, 11]:

88 **A2 (Biased Compression with contraction).** For $\delta > 0$, a compression operator is said to be
 89 $1/(1 + \delta)$ -contractive if for any $x \in \mathbb{R}^d$, we have $\mathbb{E}[\|\text{Comp}(x) - x\|] \leq (1 - 1/(1 + \delta))\|x\|$.

90 Constants ω and δ from these two assumptions are both positive, and become larger as the compression
 91 rate increases. Alternative assumptions for the biased case have been introduced in [5].

92 **Impact of unbiasedness on the compression of a single vector.**¹ To understand the interaction be-
 93 tween the number of workers K and the compression error, a simple situation is the case in which the
 94 workers use *independent and identically distributed compression operators* $(\text{Comp}_k)_{k=1}^K$ to compress
 95 the *same vector* $x \in \mathbb{R}^d$. The central node aggregates $\{\text{Comp}_k(x)\}_{k=1}^K$ into $K^{-1} \sum_{k=1}^K \text{Comp}_k(x)$.
 96 A bias-variance decomposition of the quadratic error gives:

$$\mathbb{E}[\|K^{-1} \sum_{k=1}^K \text{Comp}_k(x) - x\|^2] = \|\mathbb{E}[\text{Comp}_1(x)] - x\|^2 + K^{-1} \|\mathbb{E}[\text{Comp}_1(x)] - x\|^2.$$

97 The variance of the aggregated vector is reduced by a factor K^{-1} when averaging the messages
 98 send by the K workers, while the bias is independent of K . For example, if we use an unbiased
 99 compressor satisfying A 1, we get

$$\mathbb{E} \left[K^{-1} \sum_{k=1}^K \text{Comp}_k(x) \right] = x, \quad \mathbb{E} \left[\left\| x - K^{-1} \sum_{k=1}^K \text{Comp}_k(x) \right\|^2 \right] \leq (\omega/K) \|x\|^2, \quad (2)$$

100 while for a deterministic biased compressor, we obtain that $K^{-1} \sum_{k=1}^K \text{Comp}_k(x) = \text{Comp}_1(x)$
 101 has the same error as any of the individual compressed vector. We therefore pay particular attention
 102 to obtaining an unbiased compressor in the following.

103 2.2 StoVoQ definitions and main properties.

104 The basic idea behind VQ is to quantize a vector
 105 rather than each of its coordinates. A Vector
 106 Quantizer is a mapping $\text{VQ}(\cdot, \mathcal{C}_M) : \mathbb{R}^d \rightarrow$
 107 \mathcal{C}_M which maps $x \in \mathbb{R}^d$ to an element of a
 108 codebook \mathcal{C}_M , which is a finite subset of \mathbb{R}^d
 109 with M elements. The code of StoVoQ is pro-
 110 vided in Algorithm 1, and its crucial steps are
 111 described hereafter: we introduce the notion of
 112 (a) Voronoi quantization scheme before describ-
 113 ing more precisely (b) random codebooks, (c) whose distributions are invariant by unitary transforms.
 114 Then, (d) a method to obtain an unbiased Voronoi scheme is presented and finally (e) its asymptotic
 115 properties (as $M \rightarrow \infty$) are given.

116 (a) **Voronoi Quantization.** Voronoi quantization [23, 25], aims at selecting the closest codeword
 117 from \mathcal{C}_M , i.e.:

$$\text{VQ}(x, \mathcal{C}_M) \triangleq \underset{c \in \mathcal{C}_M}{\text{argmin}} \|x - c\|. \quad (3)$$

118 Unfortunately, for any given \mathcal{C}_M , the Voronoi quantizer is not *unbiased*: indeed it is deterministic
 119 and $\text{VQ}(x, \mathcal{C}_M) \neq x$ if $x \notin \mathcal{C}_M$. A classical approach to construct a bias-free VQ is to use the
 120 optimal ‘‘dual’’ VQ (or Delaunay quantization) [24], but this approach is numerically expensive (see
 121 Subsection 2.3). To mitigate the bias, we rather use random codebooks.

122 (b) **Random Codebook.** A key ingredient of StoVoQ is the use of a random codebook within the
 123 quantizer. We assume $\mathcal{C}_M = [C_1, \dots, C_M]$ where *the codewords* $\{C_i\}_{i=1}^M$ are i.i.d. random vectors
 124 distributed according to p , the codeword distribution pdf. We denote $\mathcal{C}_M \sim p$ and use boldface
 125 to stress that \mathcal{C}_M is random. When quantizing a sequence of vectors $\{x_t\}_{t=0}^\infty \subset \mathbb{R}^d$ we sample
 126 for each $t \in \mathbb{N}$ a **new codebook** $\mathcal{C}_{M,t} \sim p$, compute $\text{VQ}(x, \mathcal{C}_{M,t})$ and transmit the index of the
 127 corresponding codeword $i_{c,t} \in [M]$. The codebook $\mathcal{C}_{M,t}$ is **not transmitted**: the transmitter and the
 128 receiver use the **same seeds** so that the same codebooks $\mathcal{C}_{M,t}$ can be reconstructed on both sides.

¹The impact of unbiasedness for obtaining optimal convergence complexities in FL is discussed in Section 3.

Algorithm 1: StoVoQ with distribution p

Input : $x \in \mathbb{R}^d, p, M, P$, seed s

Output : Codeword index i_c , value i_r

- 1 Sample $\mathcal{C}_M \sim p$ with seed s ; /* generate
codebook with distribution p */
 - 2 $c = \text{VQ}(x, \mathcal{C}_M^p)$; /* perform Voronoi quant. */
 - 3 $i_c = \text{index of } c$; /* get index of codeword */
 - 4 $r = r_M^p(\|x\|)$; /* find radial bias in table */
 - 5 $i_r = \text{SQ}(r^{-1})$; /* quantize r on P bits */
-

129 **(c) Unitary invariant Codewords.** Denote by $U(d) = \{U, U^*U = I\}$ the set of unitary transforms
 130 over \mathbb{R}^d . We assume in the sequel that the codeword distribution p is unitary invariant, meaning that:
 131 **A3.** *The distribution of the codewords p is invariant under the unitary group, i.e. for all $U \in U(d)$,*
 132 *and any $x \in \mathbb{R}^d$, $p(Ux) = p(x)$.*

133 Examples of such distributions include isotropic Gaussian distributions ($p = \mathcal{N}(0, \sigma^2 I_d)$, $\sigma^2 > 0$)
 134 and the uniform distribution on the Sphere (which is specifically discussed in Appendix D.1). Under
 135 A 3, there exists a non-negative function p_{rad} on \mathbb{R}_+ such that, for all $x \in \mathbb{R}^d$, $p(x) = p_{\text{rad}}(\|x\|)$.

136 **(d) The quantization bias is radial.** Under A 3, we have the following crucial unitary invariance
 137 property. For $A \subset \mathbb{R}^d$, and $U \in U(d)$, we write $UA = \{Ux, x \in A\}$.

138 **Lemma 1.** *Assume A 3. For any nonnegative measurable function f , any $U \in U(d)$ and $x \in \mathbb{R}^d$,*
 139 $\mathbb{E}_{\mathcal{C}_M \sim p}[f(\text{VQ}(Ux, \mathcal{C}_M))] = \mathbb{E}_{\mathcal{C}_M \sim p}[f(U \text{VQ}(x, U \mathcal{C}_M))]$.

140 The proof is postponed to Appendix A.3. Tak-
 141 ing $f(x) = x$, the previous result implies that
 142 for any $x \in \mathbb{R}^d$ and $U \in U(d)$, it holds that
 143 $\mathbb{E}_{\mathcal{C}_M \sim p}[\text{VQ}(Ux, \mathcal{C}_M)] = U \mathbb{E}_{\mathcal{C}_M \sim p}[\text{VQ}(x, U \mathcal{C}_M)]$.
 144 A direct consequence of the elementary Lemma 3 is
 145 that the quantization error is radial:

146 **Theorem 1** (Quantization bias). *Assume A 3. Then,*
 147 *for all $M \in \mathbb{N}$, there exists a function $r_M^p : \mathbb{R}_+ \mapsto$*
 148 \mathbb{R}_+ *such that for all $x \in \mathbb{R}^d$, $\mathbb{E}_{\mathcal{C}_M \sim p}[\text{VQ}(x, \mathcal{C}_M)] =$*
 149 $r_M^p(\|x\|)x$.

150 The proof is postponed to Appendix A.4.

151 In words, the expectation of the quantized vec-
 152 tor $\text{VQ}(x, \mathcal{C}_M)$ is *colinear* to the vector x , i.e.,
 153 $\text{VQ}(x, \mathcal{C}_M)$ is **directionally unbiased**. Moreover, this radial bias only depends on $\|x\|$, M and
 154 the distribution p . This function is intractable, but it is straightforward to pre-compute it using
 155 Monte-Carlo method. We display r_M^p for $p = \mathcal{N}(0, I_d)$ in Figure 1. Consequently, we can remove
 156 the bias of $\text{VQ}(x, \mathcal{C}_M)$ by re-scaling the corresponding codeword by $1/r_M^p(\|x\|)$.

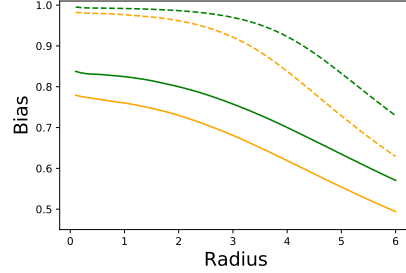


Figure 1: function r_M^p for $d = 4$ (dashed) and $d = 16$ (solid), $p = \mathcal{N}(0, I_d)$ and $M = 2^{10}$ (orange), and $M = 2^{13}$ (green).

157 We now analyze the quantization distortion for a given $x \in \mathbb{R}^d$ vector. We need to strengthen the
 158 assumption about the distribution of the codewords. Consider the following assumption

159 **A4.** (1) *there exists $\epsilon > 0$ such that $\int r^{2+\epsilon} p_{\text{rad}}(r) dr < \infty$* (2) *for some $\delta > 0$, $m_\delta =$*
 160 $\inf_{r \leq \delta} p_{\text{rad}}(r) > 0$, *and* (3) *p_{rad} is unimodal, i.e. the super level sets $\{r \in \mathbb{R}_+, p_{\text{rad}}(r) \geq t\}$,*
 161 *for $t \geq 0$ are convex subsets of \mathbb{R}_+ .*

162 A 4 is obviously satisfied if we take $p = \mathcal{N}(0, \sigma^2 I_d)$ for any $\sigma^2 > 0$.

163 **Theorem 2.** *Assume A 3-A 4. Define $C_d = \pi^{-1} \Gamma(1 + 2/d) \Gamma(1 + d/2)^{2/d}$. Then, for every $x \in \mathbb{R}^d$,*

$$\lim_{M \rightarrow \infty} M^{2/d} \mathbb{E}_{\mathcal{C}_M \sim p}[\|\text{VQ}(x, \mathcal{C}_M) - x\|^2] = C_d p_{\text{rad}}^{-2/d}(\|x\|).$$

164 The proof is postponed to Appendix C.1. Note that $C_d \approx_{d \rightarrow \infty} d/(2\pi e)$ hence C_d grows only linearly
 165 with the dimension d . We can now exploit this result to control the radial bias as a function of $\|x\|$.
 166 Since $|r_M^p(\|x\|) - 1| \leq \|x\|^{-1} \{\mathbb{E}_{\mathcal{C}_M \sim p}[\|\text{VQ}(x, \mathcal{C}_M) - x\|^2]\}^{1/2}$, Theorem 2 shows that

$$\limsup_{M \rightarrow \infty} M^{1/d} |r_M^p(\|x\|) - 1| \leq C_d^{1/2} p_{\text{rad}}^{-1/d}(\|x\|) / \|x\|.$$

167 In other words, for any $x \in \mathbb{R}^d$, the radial bias $r_M^p(\|x\|)$ approaches 1 as $M \rightarrow \infty$ with a rate
 168 $O(M^{-1/d})$. We use an a scalar quantizer SQ to transmit $1/r_M^p(\|x\|)$. Because the range of values
 169 taken by $1/r_M^p(\|x\|)$ is limited, a small number of bits P is sufficient (we typically use $P = 3$
 170 bits). The total number of transmitted bits is $\log_2(M) + \log_2(P)$. We use a random unbiased scalar
 171 quantizer (see e.g. [8, Eq. (2)]), a random mapping for $\mathbb{R} \rightarrow \mathcal{S}_P$ an ordered subset of \mathbb{R} with P
 172 elements. A scalar quantizer is said to be unbiased if $\mathbb{E}[\text{SQ}(r)] = r$ for all $r \in \mathbb{R}$. Assuming that
 173 SQ is independent of \mathcal{C}_M , we get for all $x \in \mathbb{R}^d$, $\mathbb{E}[\text{SQ}(1/r_M^p(\|x\|))] \mathbb{E}_{\mathcal{C}_M \sim p}[\text{VQ}(x, \mathcal{C}_M)] = x$. To
 174 save space, we present the details of the scalar quantization (based on nonuniform random dither)
 175 methods is presented in Appendix B.1.

176 **(e) Random vs. Optimal codebooks:** We finally motivate the choice of random codebooks and
 177 describe how to choose the codeword distribution p . For a given pdf q of the input the (*quadratic*)
 178 *distortion* is defined as:

$$\text{Dist}(q, \mathcal{C}_M) = \int_{\mathbb{R}^d} \|x - \text{VQ}(x, \mathcal{C}_M)\|^2 q(x) dx = \mathbb{E}_{X \sim q} [\|X - \text{VQ}(X, \mathcal{C}_M)\|^2]. \quad (4)$$

179 We stress that in this case the expectation is taken w.r.t. the input distribution q , the codebook
 180 being deterministic in (4). A *Voronoi optimal codebook* $\mathcal{C}_M^{q,*}$ is a minimizer of the distortion over
 181 the set of codebooks: $\text{Dist}(q, \mathcal{C}_M^{q,*}) = \min_{|\mathcal{C}_M|=M} \text{Dist}(q, \mathcal{C}_M)$. Zador's theorem [13] gives the
 182 distortion of the Voronoi optimal codebook in the limit of $M \rightarrow \infty$; see Appendix C.1 for a precise
 183 statement. Denote for $\beta \in \mathbb{R}_+$ and a function f on \mathbb{R}^d , $\|f\|_\beta = (\int |f(x)|^\beta dx)^{1/\beta}$. It is known that
 184 if $\|q\|_{d/(d+2)} < \infty$, then as $M \rightarrow \infty$, $\text{Dist}(q, \mathcal{C}_M) \approx M^{-2/d} J_d \|q\|_{d/(d+2)}$, and J_d is a universal
 185 constant J_d satisfying $J_d \approx_{d \rightarrow \infty} d/2\pi e$ (see Appendix C.2 for the exact constant).

186 Using Theorem 2, we can quantify the loss between random codebook distributed according to p and
 187 the Voronoi optimal codebook for a given input distribution q when $M \rightarrow \infty$. Define

$$C(q, p, d) = \int_{\mathbb{R}^d} p(x)^{-2/d} q(x) dx. \quad (5)$$

188 If $\|q\|_{d/(d+2)} < \infty$, using the Hölder inequality with negative exponents (see [15, p. 191] and
 189 Appendix C.3), it holds that $C(q, p, d) \geq \|q\|_{d/(d+2)}$.

190 **Theorem 3.** Assume that p satisfies A 3-A 4, $\|q\|_{d/(d+2)} < \infty$, $\int_{\mathbb{R}^d} \|x\|^{2+\delta} q(x) dx < \infty$ for some
 191 $\delta > 0$, and $C(q, p, d) < \infty$. Then,

$$\lim_{M \rightarrow \infty} \mathbb{E}_{\mathcal{C}_M \sim p} [\text{Dist}(q, \mathcal{C}_M)] / \text{Dist}(q, \mathcal{C}_M^{q,*}) = C_d J_d^{-1} C(q, p, d) \|q\|_{d/(d+2)}^{-1}. \quad (6)$$

192 with C_d defined in Theorem 2. Moreover, assume that input distribution q satisfies A 3-A 4, and set the
 193 codeword distribution $p_{q,d,*} = q^{d/(d+2)}(x) / \int q^{d/(d+2)}(x) dx$. Then, $C(q, p_{q,d,*}, d) = \|q\|_{d/(d+2)}$.

194 The proof is postponed to Appendix C.2. In words, under general assumptions, the distortion
 195 achieved by a random quantizer $\text{VQ}(\cdot, \mathcal{C}_M)$, $\mathcal{C}_M \sim p$ is rate optimal (with rate $M^{-2/d}$). If
 196 in addition q is unimodally invariant and unimodal, then a random codebook distributed accord-
 197 ing to $p_{q,d,*}$ reaches the optimal distortion bound, up to universal constants (depending only
 198 on the dimension d). Moreover, as $d \rightarrow \infty$, then $C_d J_d^{-1} \approx_{d \rightarrow \infty} 1$ and the efficiency gap van-
 199 ishes. As an illustration, assume that the input distribution is standard Gaussian $q = \mathcal{N}(0, I_d)$
 200 and set the codeword distribution to be $p_\alpha = \mathcal{N}(0, \alpha^2 I_d)$ where $\alpha^2 \in \mathbb{R}_+^*$. If $\alpha^2 d > 2$, then
 201 $C(\mathcal{N}(0, I_d), \mathcal{N}(0, \alpha^2 I_d), d) = 2\pi\alpha^2 \{\alpha^2 d / (\alpha^2 d - 2)\}^{d/2}$ and $\|\mathcal{N}(0, I_d)\|^{(2+d)/2} = (2\pi)(1 +$
 202 $2/d)^{1+2/d}$. The function $\alpha \rightarrow C(\mathcal{N}(0, I_d), \mathcal{N}(0, \alpha^2 I_d), d)$ has a unique minimum at $\alpha_d^2 = 1 + 2/d$
 203 for which $C(\mathcal{N}(0, I_d), \mathcal{N}(0, \alpha_d^2 I_d), d) = \|\mathcal{N}(0, I_d)\|^{(2+d)/2}$ showing that a random codebook sam-
 204 pled from $\mathcal{N}(0, \alpha_d^2 I_d)$ is optimal. It is interesting to note that the variance of the codeword distribution
 205 should be $(1 + 2/d)$ larger than the variance of the input distribution $\mathcal{N}(0, I_d)$.

206 2.3 Related works

207 We compare StVoQ with competing (random) compressors; additional details are given App. A.1.

208 **QSGD.** Alistarh et al. [2] compresses each coordinate of the scaled vector $x/\|x\|$ on $s+1$ codewords.
 209 QSGD is a scalar quantizer which requires $\mathcal{O}(\sqrt{d} \log_2(d))$ bits in its highest compression setting
 210 ($s = 1$, only two possible levels for each coordinate). The vector norm is transmitted with full
 211 precision $\|x\|$ (16 or 32 bits). This is in general substantially higher than the number of bits used by
 212 VQ methods. In deep learning problems, it reduces the communication cost by a factor of 4 to 7 [2,
 213 Sec. 5].

214 **Top-H/Rand H.** Achieving higher compression rates is possible through *sparsification* operators, that
 215 only transmit a few coordinates. The most popular schemes are Top- H and Rand- H compressors,
 216 that respectively map the vector to either its H largest coordinates, or a random subset of cardinality
 217 H , rescaled by d/H to ensure unbiasedness. Top- H is a biased operator, and the performance of
 218 Rand- H are poor on deep learning tasks [5, Figures 4 and 5].

Table 1: Per iteration communication complexity of most frequently used algorithms in dimension d . Constants H and M respectively correspond to a number of coordinates to be transmitted and a number of codewords, they are chosen by the user.

#bits	Uncomp.		Scalar Quantization			Vector Quantization				
	SGD	Sign	QSGD $s \geq 1$	Top- H	Rand- H	Polytope [10]	HSQ-span [8]	HSQ-greed [8]	StoVoQ	DoStoVoQ
	$32d$	d	$32 + s\sqrt{d}\log_2(d)$	$32H$	$32H$	$\log_2(2d)$	$\log_2(M)$	$\log_2(M)$	$\log_2(M)$	$\log_2(M)$
Unbiased	-	-	\checkmark	-	\checkmark	\checkmark	\checkmark	-	\checkmark	\checkmark (Th.4)
A.1 ($\omega + 1$)	-	-	\sqrt{d}/s	-	d/H	d	d	-	-	$O(M^{-2/d})$ (Th.4)
A.2 ($\delta + 1$)	-	-	-	d/H	-	-	-	$M/\sigma_{\min}(C)$	-	-

219 **HyperSphere Quantization (HSQ)**. HSQ was introduced by Dai et al. [8]. Two versions are consid-
 220 ered: (1) a - greedy- Voronoi VQ referred to as HSQ-greed in Table 1, which is biased, and for which
 221 the theoretical guarantee provided in the paper (in their Lemma 3 and Theorem 3, which corresponds
 222 to a variant of A 2 and the subsequent convergence rate) *worsens* as M increases, making it mostly
 223 vacuous; (2) an unbiased version VQ (HSQ-span), which uses a minimum-norm decomposition of
 224 $x \in \text{Span}(\mathcal{C}_M)$ the linear subspace generated by the codewords - this version suffers from a large
 225 variance (see Table 2) and potentially an ill-conditioning. Moreover, the performance of HSQ-span
 226 does not improve with M .

227 StoVoQ builds on HSQ-greed, that achieves high compression factors (up to 60-100 to obtain close
 228 to SOTA performance on CIFAR10), while preserving a good flexibility w.r.t. the compression
 229 level. StoVoQ approach allows to remove its inherent bias and provide a much stronger convergence
 230 analysis: **our approach is the first vector quantization scheme to provably benefit from an**
 231 **increasing number of elements in the codebook M** (and obviously benefits from the number of
 232 workers K , as it is unbiased).

233 **Dual Quantization and Cross-polytope**. An approach to constructing unbiased VQ is to use
 234 the dual VQ, also referred to as Delaunay Quantization (DQ); see [24]. DQ is unbiased for any
 235 $x \in \text{ConvHull}(\mathcal{C}_M)$, the convex hull of \mathcal{C}_M . DQ requires to compute the barycentric coordinates
 236 for $x \in \text{ConvHull}(\mathcal{C}_M)$, that is to solve $(\lambda_1^x, \dots, \lambda_M^x) = \text{argmin}_{\lambda_1, \dots, \lambda_M} \|x - \sum_{i=1}^M \lambda_i c_i\|^2$, under
 237 the constraints $\lambda_i \geq 0$, $\sum_{i=1}^M \lambda_i = 1$. The quantizer is obtained by drawing a codeword c_i with
 238 probability $[\lambda_1^x, \dots, \lambda_M^x]$. Computing the barycentric coordinates is in general very demanding
 239 unless \mathcal{C}_M has a very simple structure (see Appendix B for details). The Cross-Polytope
 240 method Gandikota et al. [10] is a simple instance of DQ, with a codebook $\mathcal{C}_{2d}^{\text{CP}}$ composed of the
 241 $2d$ canonical vectors $\{\pm \sqrt{d}e_i = \pm(0, \dots, 0, \sqrt{d}, 0 \dots 0), i \in [d]\}$, that relies on the inclusion
 242 $B_2(0; 1) \subset B_1(0; \sqrt{d}) = \text{ConvHull}(\mathcal{C}_{2d}^{\text{CP}})$. The barycentric decomposition can then easily be
 243 computed. Unfortunately, this method suffers from a large variance, as the quantization error
 244 $\|\text{VQ}^{\text{CP}}(x, \mathcal{C}_M) - x\|$ of *any* x is *lower bounded* by $\sqrt{d} - 1$, which means the error has the same
 245 quadratic error than the Rand-1 compressor.

246 Table 1 summarizes the number of bits required to exchange the compressed value of a vector $x \in \mathbb{R}^d$
 247 for the compression methods considered in this Section, as well as the assumptions they satisfy.

248 **Numerical comparisons:** In Table 2, we compare the distortions achieved by the compression
 249 methods given in Table 1 for a communication budget of 16 bits for $d = 16$ and assuming that the
 250 input distribution is $q = \mathcal{N}(0, I_d)$. The compression factor is 32 (assuming 32 bits floating point
 251 per coordinate). Such a compression rate is out of reach for QSGD, that requires, even for $s = 1$ at
 252 least $\sqrt{d}\log(d) + R$ bits, where R is the number of bits to encode the norm (32 in [2]). For QSGD we
 253 have quantized the norm (using an uniform quantizer) on 3 bits and obtained an averaged distortion
 254 of 36.10 (for $K = 1$) and 1.82 for ($K = 20$) - the total number of bits is 19-. We use $H = 2$ for
 255 Top- H and Rand- H and use a scalar quantizer with 8 bits. For HSQ, we use 6 bits for the norm,
 256 using the unbiased uniform quantizer given in [8] and a Voronoi optimal codebook for the uniform
 257 distribution on the unit-sphere with $M = 2^{10}$ codewords. For StoVoQ we use a random codebook
 258 with $M = 2^{13}$ codewords; the codewords are sampled from a $\mathcal{N}(0, (1 + 2/d)I_d)$, and 3 bits are
 259 allocated for the scalar quantization of $1/r_M^p$ (the inverse of the radial bias). Finally, we average the
 260 result of 2 independent compressions for Polytope (following the replication technique described in
 261 [10]). We use $n = 10^4$ vectors, and report in Table 2 the distortion and sample variance. For StoVoQ
 262 with $K = 20$, the codebooks of the different workers are independent.

Table 2: Distortion for Gaussian inputs, for a fixed budget of 16 bits with $d = 16$.

Method	Sign [4]	Top-2	Rand-2	Polytope [10]	HSQ-span [8]	HSQ-greed [8]	StoVoQ
# Bits (obj=16)	16	2×8	2×8	$\log_2(2 \times 16) \times 2 + 6$	$\log_2(2^{10}) + 6$	$\log_2(2^{10}) + 6$	$\log_2(2^{13}) + 3$
Unbiased			✓	✓	✓		✓
$K = 1$	6.21 (0.02)	8.40 (0.04)	102.8 (0.9)	113.9 (0.6)	146.9 (0.6)	9.03 (0.04)	6.97 (0.02) :
$K = 20$	6.26 (0.02)	8.76 (0.04)	5.40 (0.04)	5.98 (0.03)	7.58 (0.04)	9.10 (0.04)	0.838 (0.005)

263 3 DoStoVoQ algorithm

264 We illustrate how the StoVoQ compression scheme can be implemented in FL. To avoid cumbersome
 265 technical details, we focus here on the Federated-SGD algorithm. At iteration $t + 1$, each worker
 266 computes a stochastic gradient $g_{k,t+1}$ of the loss f_k at the current model θ_t , compresses it into
 267 $\hat{g}_{k,t+1} = \text{Comp}(g_{k,t+1})$ and send it to the central server, that performs the update step $\theta_t =$
 268 $\theta_{t-1} - \gamma_t / K \sum_{k=1}^K \hat{g}_{k,t}$. The code of the resulting algorithm, DoStoVoQ-SGD, is given in Algorithm 2.
 269 At iteration $t + 1$, the crucial steps are:

- 270 1. Worker $k \in [K]$ computes the norm $\|g_{k,t+1}\|$ of the $D \times 1$ gradient $g_{k,t+1}$ and then splits
 271 the scaled gradient $g_{k,t+1} \times \sqrt{D} / \|g_{k,t+1}\|$ into L -buckets of size d : $g_{k,t+1} \times \sqrt{D} / \|g_{k,t+1}\| =$
 272 $[b_{k,t+1}^1, \dots, b_{k,t+1}^L]$. The norm $\|g_{k,t+1}\|$ is transmitted to the central node using a high-resolution
 273 scalar quantizer (or without quantization).
- 274 2. Each worker quantizes the buckets $\{b_{k,t+1}^1, \dots, b_{k,t+1}^L\}$ using StoVoQ. **Independent** codebooks
 275 $\{\mathcal{C}_{M,k,t+1}\}_{k \in [K]}$ are used to ensure that the quantizers remain conditionally independent (see
 276 below for a precise statement). The double stochasticity (each worker uses random codebooks,
 277 which are independent between workers and across iterations) motivates the name DoStoVoQ. At
 278 iteration t , the same codebook is used for all buckets of worker k . Formally, for $\ell \in [L]$ we apply
 279 (in parallel) $\text{StoVoQ}(b_{k,t+1}^\ell, p, M, P, s_{k,t+1})$, with a sequence of different seeds $(s_{k,t+1})_{k \in [K], t \geq 0}$.
 280 This sequence is shared between the workers and the central node at initialization.
- 281 3. The central node computes $(\hat{g}_{k,t+1})_{k \in [K]}$ from all messages received, performs the update on
 282 $(\theta_t)_{t \geq 0}$, and broadcasts θ_{t+1} to the workers.

283 These steps would similarly allow to incorporate StoVoQ within any of the advanced FL algo-
 284 rithms, and Theorem 4 is the crucial assumption to derive the convergence rates, as described in
 285 Section 2. Natural extensions to DoStoVoQ-Fed-Avg, DoStoVoQ-DIANA and DoStoVoQ-VR-DIANA
 286 are provided in Appendix D.2.

287 **Bias and variance of the com-**
 288 **pressed gradient with K workers.**
 289 Consider the two filtrations $(\mathcal{F}_t)_{t \geq 0}$
 290 and $(\mathcal{G}_t)_{t \geq 0}$ defined recursively as fol-
 291 lows $\mathcal{F}_0 = \sigma(\emptyset)$ and for $t \geq 0$,
 292 $\mathcal{G}_{t+1} = \mathcal{F}_t \vee \sigma(\{g_{k,t+1}, k \in [K]\})$
 293 and $\mathcal{F}_{t+1} = \mathcal{G}_{t+1} \vee \sigma(\{\hat{g}_{k,t+1}, k \in$
 294 $[K]\})$. With these notations, for any
 295 $t \geq 0$, θ_t is \mathcal{F}_t -measurable.

296 **Theorem 4.** *At any iteration $t +$
 297 1 in DoStoVoQ, the K compressed
 298 stochastic gradients $(\hat{g}_{k,t+1})_{k \in [K]}$
 299 are (i) independent conditionally
 300 to \mathcal{G}_{t+1} (ii) conditionally unbiased,
 301 i.e., for all $k \in [K]$, we have
 302 $\mathbb{E}[\hat{g}_{k,t+1} | \mathcal{G}_{t+1}] = g_{k,t+1}$, (iii) sat-
 303 isfy the relatively bounded error con-
 304 dition of A 1, i.e. there exists a con-
 305 stant ω_M such that, for all $k \in [K]$: $\mathbb{E}[\|\hat{g}_{k,t+1} - g_{k,t+1}\|^2 | \mathcal{G}_{t+1}] \leq \omega_M \|g_{k,t+1}\|^2$.*

306 *Moreover, ω_M decreases with the number of codewords M and the P , as $\omega_M = O(M^{-2/d}) + O(2^{-P})$
 307 [the dependence on p, d , and D is made explicit in the proof].*

Algorithm 2: DoStoVoQ-SGD over T iterations

Input : T nb of steps, $(\gamma_t)_{t \geq 0}$ LR, θ_0, p, M, P ;
Output : $(\theta_t)_{t \geq 0}$

- 1 **for** $t = 1, \dots, T$ **do**
- 2 w_0 sends θ_{t-1} and different seeds $s_{k,t}$ to each w_k ;
- 3 **for** $k = 1, \dots, K$ **do**
- 4 Compute local gradient $g_{k,t}$ at θ_{t-1} ;
- 5 Split $g_{k,t} \times \sqrt{D} / \|g_{k,t}\|$ on $[b_{k,t}^1, \dots, b_{k,t}^L]$;
- 6 **for** $\ell = 1, \dots, L$ (in parallel) **do**
- 7 $(\mathbf{i}_c^{t,k,\ell}, \mathbf{i}_r^{t,k,\ell}) = \text{StoVoQ}(b_{k,t}^\ell, p, M, P, s_{k,t})$
- 8 **end**
- 9 Send $(\|g_{k,t}\|, (\mathbf{i}_c^{t,k,\ell}, \mathbf{i}_r^{t,k,\ell})_{\ell \in [L]})$ to w_0 ;
- 10 **end**
- 11 Reconstruct $(\hat{g}_{k,t})_{k \in [K]}$;
- 12 Update: $\theta_t = \theta_{t-1} - \gamma_t \frac{1}{K} \sum_{k=1}^K \hat{g}_{k,t}$;
- 13 **end**

308 The first statement stems from the fact that each bucket is quantized using StoVoQ which is unbiased.
 309 The second statement is more challenging; proof is postponed to Appendix A.6. We stress that this
 310 result differs from Theorem 2, which corresponds to the distortion of a source with distribution q .

311 **Convergence results.** Theorem 4 proves that our compression method satisfies the assumptions
 312 needed to obtain fast convergence rate, for DoStoVoQ-SGD, and for its variants DoStoVoQ-(VR)-
 313 DIANA. Consider a Smooth and Strongly Convex (SSC) function $F = \sum_{k=1}^K f_k$, with condition
 314 number $\kappa > 1$. We measure the complexity of the algorithm by the number of iterations t required
 315 to obtain a model θ_t such that $\mathbb{E}[F(\theta_t)] - \min_{\mathbb{R}^D} F \leq \epsilon$. The result of VR-DIANA [16], which
 316 provides a complexity of $O_{\kappa \rightarrow \infty}(\kappa(1 + \omega_M/K) \log(\epsilon^{-1}))$ [16, Corollary 2], applies to DoStoVoQ-
 317 VR-DIANA.

318 Convergence rates for DoStoVoQ-DIANA (without VR), and on non-convex optimization problems
 319 can be obtained from Horváth et al. [16, Corollary 1,3,4]. As in the strongly-convex case, complexities
 320 increase by a factor depending on $(1 + \omega_M/K)$ w.r.t. uncompressed algorithm. Intuitively, *the impact*
 321 *on the optimization complexity of a high compression is mitigated by the number of workers*, which
 322 supports the use of independent and unbiased compressors when the number of workers is large and
 323 high compression factors are required.

324 Indeed, these complexities can be compared to: (1) the one of *uncompressed* variance reduced
 325 distributed methods [9] that achieve a complexity of $O_{\kappa \rightarrow \infty}(\kappa \log(\epsilon^{-1}))$ (in the SSC case); (2) the
 326 complexity for biased compression operators satisfying A 2, Beznosikov et al. [5, Theorem 13] that
 327 obtain $O_{\kappa \rightarrow \infty}(\kappa(1 + \delta) \log(\epsilon^{-1}))$ for compressed GD (independently of the number of workers);
 328 (3) the complexities of compressed SGD methods with *error feedback* in [11]², that also have no
 329 dependency on the number of workers. **Overall, the unbiased character is crucial to mitigate the**
 330 **variance increase resulting from high compression rates.**

331 4 Numerical experiments

332 4.1 Least Squares Regression (LSR)

333 We consider a least-squares problem with $n =$
 334 2^{14} samples, a bucket size $d = 16$, $D = 2^9$, and
 335 $K = 32$ workers; each worker has access to a
 336 subset $m = 2^{11}$ samples (picked with replace-
 337 ment) to introduce a dependency in the data used
 338 by the workers. For $i \in [n]$, we assume $X_i \sim$
 339 $\mathcal{N}(0, I_D)$ and $Y_i \sim \mathcal{N}(X_i^\top \omega_*, 1)$ where $\omega_* \in$
 340 \mathbb{R}^D . We solve $\inf_{\omega \in \mathbb{R}^D} \sum_{i=1}^n \|Y_i - X_i^\top \omega\|^2$ via
 341 a gradient descent with step size $1/\alpha L$ where
 342 α is fine-tuned for each quantization method
 343 and $L \approx 2n$ is the smoothness constant. We
 344 use DoStoVoQ with $M = 2^{13}$ codewords sam-
 345 pled from $\mathcal{N}(0, (1 + 2/d) I_d)$ for DoStoVoQ and
 346 $M = 2^{10}$ on the unit Sphere for HSQ s.t. the
 347 number of bits transmitted at each round by the
 348 worker is set to 16 (see Table 2). Figure 2 reports
 349 the excess-log of the train loss over $T = 10$ iterations, for a standard GD. DoStoVoQ outperforms
 350 HSQ-greed: indeed the linear convergence rate of distributed GD is faster for an unbiased compressor
 351 than for the biased approach.

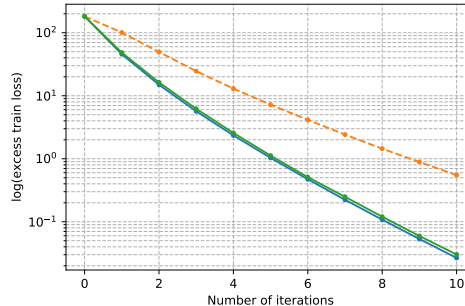


Figure 2: Comparison between GD (blue), HSQ-greed (orange) and DoStoVoQ (green), on a LSR problem in dimension $D = 2^9$.

352 4.2 Applications to Deep Neural Networks training

353 **Setting.** We now describe our experimental framework for training two standard models of Deep
 354 Neural Networks: a VGG-16 [31] and a ResNet-18 [14]. We follow the standard procedure of training
 355 those models both on CIFAR-10 and ImageNet; the hyper-parameters are fine-tuned to optimize the
 356 accuracy *without quantization*. We do not compress the affine constant part of the affine convolutional

²authors provide complexities for 10 algorithms in Table 1, with Error Feedback and under A 2.

Table 3: Average accuracy over 5 experiments, after 100 epochs on CIFAR-10.

Algorithm	SGD	QSGD	QSGD	QSGD	HSQ	HSQ	Dos.	Dos.
		2 bits	4 bits	8 bits	$d = 16$	$d = 8$	$d = 16$	$d = 8$
Raw bits per bucket	$32d$	$\sqrt{d} \log(d)$			$\log(d)$			
Effective Compression factor	1	~ 13	~ 8	~ 4	34	17	38	20
$K = 1$ worker	91.9	91.7	92.1	91.9	92.0	92.0	92.0	92.1
$K = 8$ worker	92.0	91.8	91.8	92.0	91.8	92.0	91.8	92.1

Table 4: Distortion for on a subset \mathcal{G} of the gradients of a layer of CIFAR-10, for a fixed budget of 16 bits with $d = 16$.

Method	Top-2	Rand-2	Polytope [10]	HSQ-span [8]	HSQ-greed [8]	DoStoVoQ
# Bits (obj =16)	2×8	2×8	$\log_2(2 \times 16) \times 2 + 6$	$\log_2(2^{10}) + 6$	$\log_2(2^{10}) + 6$	$\log_2(2^{13}) + 3$
Unbiased		✓	✓	✓		✓
$K = 1$	0.0022	0.025	0.028	0.034	0.0021	0.0026

357 layers and batch normalization layers. We apply independent DoStoVoQ on batches of 32 buckets of
 358 size $d = 16$ (i.e. we transmit a high-resolution norm for $D = 32 \cdot 16 = 512$ coefficients).

359 **CIFAR-10.** We use the implementation of HSQ [8]: the batch size is 256 for CIFAR-10, the total
 360 number of epochs is 100, the initial learning rate is 0.1, which is divided by 10 and 50 at epochs
 361 51 and 71. We report the accuracy of DoStoVoQ, QSGD, and HSQ-greed in table 4. By design, the
 362 compression factor of Q-SGD for $d = 16$ is 13, which is significantly less than HSQ or DoStoVoQ.
 363 Both HSQ and DoStoVoQ perform similarly and the accuracy gap between the two methods are under
 364 the sample variance (computed over 5 seed and about 0.2). In Table 4 we report the distortion of
 365 a random subset of gradients $\mathcal{G} = \{g_t, t \in [|\mathcal{G}|]\}$ (with $|\mathcal{G}| = 10^2$, $d = 16$, $D = 2^5 \times d$) obtained
 366 from a given layer of a VGG on CIFAR-10, i.e.: $|\mathcal{G}|^{-1} \sum_{g_t \in \mathcal{G}} \left\| K^{-1} \sum_{k=1}^K (g_{k,t} - \hat{g}_{k,t}) \right\|^2$, where
 367 $(\hat{g}_{k,t})_{k \in [K]}$ correspond to k independent workers compressing their own gradient $g_{k,t}$. The choice
 368 of the layer does not affect significantly the results. Even with the actual gradient distribution,
 369 DoStoVoQ outperforms for a given compression factor each unbiased method. This is on pair
 370 with the observation that the gradients of a Deep Neural Network are approximately Gaussian
 371 distributed [3, 36, 4]. Additional experiments can be found in the Appendix.

372 **ImageNet.** For ImageNet, we use different bucket sizes, the standard batch size of 256, and only
 373 $K = 1$ worker for energy savings (recall Imagenet training last about 1 day for a single worker on
 374 academic hardware). An initial learning rate of 0.1 is divided by 10 at epoch 30 and 60, while the
 375 model is trained for 90 epochs. A ResNet here obtains 69.9%, and with a compression factor of 8,
 376 the performance drops by 2.5%. Using $d = 16$, we reach a compression factor of 38, while the Top-1
 377 accuracy drops by only 4.8%: this is a substantially higher compression rate than the concurrent work
 378 QSGD on the ImageNet dataset.

379 **Computational impact.** In the case of deep Neural Networks, our training procedure requires
 380 neither a substantial modifications of standard pipelines, nor a modification of the hyper-parameters
 381 which allows to save computational resources. Green Algorithm ([20]) shows that this work
 382 generated around 15kg of CO2, and require 400 kWh. A typical experiment lasted few hours on
 383 CIFAR-10 and about 3 days on ImageNet, which is in the standard range for this type of prototypical
 384 codes. This work could have future impact on FL, to reduce their electrical consumption.

385 **Broader impact.** Federated learning enables multiple actors to build a common model without
 386 data sharing, hence respecting privacy. However classic FL methods consume an important amount
 387 of energy in transmitting information. Our method DoStoVoQ can be adapted to any FL framework
 388 while enabling important bandwidth savings. These savings highly counterbalance the computational
 389 impact of our experiments.

References

- 390
- 391 [1] Naman Agarwal, Ananda Theertha Suresh, Felix Xinnan X Yu, Sanjiv Kumar, and Brendan
392 McMahan. cpSGD: Communication-efficient and differentially-private distributed SGD. In
393 S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors,
394 *Advances in Neural Information Processing Systems 31*, pages 7564–7575. Curran Associates,
395 Inc., 2018.
- 396 [2] Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. QSGD:
397 Communication-Efficient SGD via Gradient Quantization and Encoding. *Advances in Neural*
398 *Information Processing Systems*, 30:1709–1720, 2017.
- 399 [3] Ron Banner, Itay Hubara, Elad Hoffer, and Daniel Soudry. Scalable methods for 8-bit training
400 of neural networks. In *Proceedings of the 32nd International Conference on Neural Information*
401 *Processing Systems*, pages 5151–5159, 2018.
- 402 [4] Jeremy Bernstein, Yu-Xiang Wang, Kamyar Aizzadenesheli, and Animashree Anandkumar.
403 signsgd: Compressed optimisation for non-convex problems. In *International Conference on*
404 *Machine Learning*, pages 560–569. PMLR, 2018.
- 405 [5] Aleksandr Beznosikov, Samuel Horváth, Peter Richtárik, and Mher Safaryan. On Biased
406 Compression for Distributed Learning. *arXiv:2002.12410 [cs, math, stat]*, February 2020.
407 arXiv: 2002.12410.
- 408 [6] Léon Bottou. On-line learning and stochastic approximations. 1999. doi: 10.1017/
409 CBO9780511569920.003.
- 410 [7] Léon Bottou. Large-Scale Machine Learning with Stochastic Gradient Descent. In Yves
411 Lechevallier and Gilbert Saporta, editors, *Proceedings of COMPSTAT’2010*, pages 177–
412 186, Heidelberg, 2010. Physica-Verlag HD. ISBN 978-3-7908-2604-3. doi: 10.1007/
413 978-3-7908-2604-3_16.
- 414 [8] Xinyan Dai, Xiao Yan, Kaiwen Zhou, Han Yang, Kelvin KW Ng, James Cheng, and Yu Fan.
415 Hyper-sphere quantization: Communication-efficient sgd for federated learning. *arXiv preprint*
416 *arXiv:1911.04655*, 2019.
- 417 [9] Aaron Defazio, Francis R Bach, and Simon Lacoste-Julien. Saga: A fast incremental gradient
418 method with support for non-strongly convex composite objectives. In *NIPS*, 2014.
- 419 [10] Venkata Gandikota, Daniel Kane, Raj Kumar Maity, and Arya Mazumdar. vqsgd: Vector
420 quantized stochastic gradient descent. In *International Conference on Artificial Intelligence*
421 *and Statistics*, pages 2197–2205. PMLR, 2021.
- 422 [11] Eduard Gorbunov, Dmitry Kovalev, Dmitry Makarenko, and Peter Richtárik. Linearly Converging
423 Error Compensated SGD. *arXiv:2010.12292 [cs, math]*, October 2020. arXiv: 2010.12292.
- 424 [12] Eduard Gorbunov, Konstantin Burlachenko, Zhize Li, and Peter Richtárik. Marina: Faster
425 non-convex distributed learning with compression. *arXiv preprint arXiv:2102.07845*, 2021.
- 426 [13] Siegfried Graf and Harald Luschgy. *Foundations of quantization for probability distributions*.
427 Springer, 2007.
- 428 [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for im-
429 age recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*
430 *Recognition (CVPR)*, June 2016.
- 431 [15] Edwin Hewitt and Karl Stromberg. *Real and abstract analysis: a modern treatment of the*
432 *theory of functions of a real variable*. Springer-Verlag, 2013.
- 433 [16] Samuel Horváth, Dmitry Kovalev, Konstantin Mishchenko, Sebastian Stich, and Peter
434 Richtárik. Stochastic Distributed Learning with Gradient Quantization and Variance Reduction.
435 *arXiv:1904.05115 [math]*, April 2019. arXiv: 1904.05115.

- 436 [17] Peter Kairouz, H. Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Ar-
437 jun Nitin Bhagoji, Keith Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings,
438 Rafael G. L. D’Oliveira, Salim El Rouayheb, David Evans, Josh Gardner, Zachary Garrett,
439 Adrià Gascón, Badih Ghazi, Phillip B. Gibbons, Marco Gruteser, Zaid Harchaoui, Chaoyang
440 He, Lie He, Zhouyuan Huo, Ben Hutchinson, Justin Hsu, Martin Jaggi, Tara Javidi, Gauri
441 Joshi, Mikhail Khodak, Jakub Konečný, Aleksandra Korolova, Farinaz Koushanfar, Sanmi
442 Koyejo, Tancrède Lepoint, Yang Liu, Prateek Mittal, Mehryar Mohri, Richard Nock, Ayfer
443 Özgür, Rasmus Pagh, Mariana Raykova, Hang Qi, Daniel Ramage, Ramesh Raskar, Dawn
444 Song, Weikang Song, Sebastian U. Stich, Ziteng Sun, Ananda Theertha Suresh, Florian Tramèr,
445 Praneeth Vepakomma, Jianyu Wang, Li Xiong, Zheng Xu, Qiang Yang, Felix X. Yu, Han Yu,
446 and Sen Zhao. Advances and Open Problems in Federated Learning. *arXiv:1912.04977 [cs,*
447 *stat]*, December 2019. arXiv: 1912.04977.
- 448 [18] Jakub Konečný, H. Brendan McMahan, Daniel Ramage, and Peter Richtárik. Federated
449 Optimization: Distributed Machine Learning for On-Device Intelligence. *arXiv:1610.02527*
450 *[cs]*, October 2016. arXiv: 1610.02527.
- 451 [19] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
452 2009.
- 453 [20] Loïc Lannelongue, Jason Grealey, and Michael Inouye. Green Algorithms: Quantifying the
454 carbon emissions of computation. *arXiv:2007.07610 [cs]*, October 2020. arXiv: 2007.07610.
- 455 [21] Zhize Li, Dmitry Kovalev, Xun Qian, and Peter Richtarik. Acceleration for Compressed
456 Gradient Descent in Distributed and Federated Optimization. In *International Conference on*
457 *Machine Learning*, pages 5895–5904. PMLR, November 2020. ISSN: 2640-3498.
- 458 [22] Konstantin Mishchenko, Eduard Gorbunov, Martin Takáč, and Peter Richtárik. Distributed
459 Learning with Compressed Gradient Differences. *arXiv:1901.09269 [cs, math, stat]*, June 2019.
460 arXiv: 1901.09269.
- 461 [23] Gilles Pagès and Jacques Printems. Optimal quadratic quantization for numerics: the gaussian
462 case. *Monte Carlo methods and applications*, 9(2):135–165, 2003.
- 463 [24] Gilles Pagès and Benedikt Wilbertz. Sharp rate for the dual quantization problem. In *Séminaire*
464 *de Probabilités XLIX*, volume 2215 of *Lecture Notes in Math.*, pages 405–454. Springer, Cham,
465 2018.
- 466 [25] Gilles Pagès and Benedikt Wilbertz. Sharp rate for the dual quantization problem. In *Séminaire*
467 *de Probabilités XLIX*, pages 405–454. Springer, 2018.
- 468 [26] Constantin Philippenko and Aymeric Dieuleveut. Artemis: tight convergence guarantees for
469 bidirectional compression in Federated Learning. *arXiv:2006.14591 [cs, stat]*, November 2020.
470 arXiv: 2006.14591.
- 471 [27] Ali Ramezani-Kebrya, Fartash Faghri, and Daniel M Roy. Nuqsgd: Improved communication
472 efficiency for data-parallel sgd via nonuniform quantization. *arXiv preprint arXiv:1908.06077*,
473 2019.
- 474 [28] Herbert Robbins and Sutton Monro. A Stochastic Approximation Method. *Annals of Math-*
475 *ematical Statistics*, 22(3):400–407, September 1951. ISSN 0003-4851, 2168-8990. doi:
476 10.1214/aoms/1177729586. Number: 3 Publisher: Institute of Mathematical Statistics.
- 477 [29] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
478 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual
479 recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- 480 [30] F. Seide, H. Fu, Jasha Droppo, G. Li, and D. Yu. 1-bit stochastic gradient descent and its
481 application to data-parallel distributed training of speech DNNs. pages 1058–1062, January
482 2014.
- 483 [31] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale
484 image recognition. *arXiv preprint arXiv:1409.1556*, 2014.

- 485 [32] Sebastian U. Stich and Sai Praneeth Karimireddy. The Error-Feedback Framework: Better
 486 Rates for SGD with Delayed Gradients and Compressed Communication. *arXiv:1909.05350*
 487 [*cs, math, stat*], September 2019. arXiv: 1909.05350.
- 488 [33] Sebastian U Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified SGD with Memory.
 489 In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors,
 490 *Advances in Neural Information Processing Systems 31*, pages 4447–4458. Curran Associates,
 491 Inc., 2018.
- 492 [34] Thijs Vogels, Sai Praneeth Karimireddy, and Martin Jaggi. Powersgd: Practical low-rank
 493 gradient compression for distributed optimization. *Advances in Neural Information Processing*
 494 *Systems*, 32:14259–14268, 2019.
- 495 [35] Jianqiao Wangni, Jialei Wang, Ji Liu, and Tong Zhang. Gradient Sparsification for
 496 Communication-Efficient Distributed Optimization. *Advances in Neural Information Pro-*
 497 *cessing Systems*, 31:1299–1309, 2018.
- 498 [36] An Xu, Zhouyuan Huo, and Heng Huang. Optimal gradient quantization condition for
 499 communication-efficient distributed training. *arXiv preprint arXiv:2002.11082*, 2020.
- 500 [37] Hantian Zhang, Jerry Li, Kaan Kara, Dan Alistarh, Ji Liu, and Ce Zhang. The zipml framework
 501 for training models with end-to-end low precision: The cans, the cannots, and a little bit of deep
 502 learning. *arXiv preprint arXiv:1611.05402*, 2016.

503 Checklist

- 504 1. For all authors...
- 505 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 506 contributions and scope? [Yes] See Section 2 for quantization and Section 4 for
 507 associated experiments.
- 508 (b) Did you describe the limitations of your work? [Yes] See broader impact and Appendix.
- 509 (c) Did you discuss any potential negative societal impacts of your work? [Yes] Detailed
 510 experiments carbon footprint can be find in Section 4.
- 511 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 512 them? [Yes]
- 513 2. If you are including theoretical results...
- 514 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 515 (b) Did you include complete proofs of all theoretical results? [Yes] Also see Appendix in
 516 Supplemental Material.
- 517 3. If you ran experiments...
- 518 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 519 mental results (either in the supplemental material or as a URL)? [Yes] Code available
 520 in Supplementary Material.
- 521 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 522 were chosen)? [Yes] See Section 4.
- 523 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 524 ments multiple times)? [Yes] In particular Table 2 presents standard deviations, and
 525 variances of NN model accuracies from Section 4 can be found in Appendix.
- 526 (d) Did you include the total amount of compute and the type of resources used (e.g.,
 527 type of GPUs, internal cluster, or cloud provider)? [Yes] See Section 4 for further
 528 references.
- 529 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 530 (a) If your work uses existing assets, did you cite the creators? [Yes] As mentioned in
 531 Section 4, code is partly inspired from [8].
- 532 (b) Did you mention the license of the assets? [Yes] Only open source and/or Academic
 533 assets are used.

- 534 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
535 Radial biases already computed available in Supplemental Material.
- 536 (d) Did you discuss whether and how consent was obtained from people whose data
537 you're using/curating? [N/A] Use of publicly available data (CIFAR10 [19] and
538 Imagenet [29]).
- 539 (e) Did you discuss whether the data you are using/curating contains personally identifiable
540 information or offensive content? [N/A]
- 541 5. If you used crowdsourcing or conducted research with human subjects...
- 542 (a) Did you include the full text of instructions given to participants and screenshots, if
543 applicable? [N/A]
- 544 (b) Did you describe any potential participant risks, with links to Institutional Review
545 Board (IRB) approvals, if applicable? [N/A]
- 546 (c) Did you include the estimated hourly wage paid to participants and the total amount
547 spent on participant compensation? [N/A]