
DoStoVoQ: Doubly Stochastic Voronoi Vector Quantization SGD for Federated Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

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

1 Introduction

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

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

where each function $f_k : \mathbb{R}^D \rightarrow \mathbb{R}$ represents the empirical risk on worker $k \in [K]$ (where $[K] = \{1, \dots, K\}$) and D is the ambient dimension of our problem. Each worker potentially holds a fraction of the data, and can share information with a central server, which progressively aggregates and updates the model accordingly [20, 19].

Stochastic gradient algorithms [32] are particularly well suited in the *large scale learning* setting [6, 7]. The methods can easily be adapted to the distributed (and more generally federated) learning framework; see [19] and the references therein. For synchronous distributed Stochastic Gradient Descent, at every iteration, given the current parameter θ_t , each worker computes an unbiased estimate $g_{k,t+1}(\theta_t)$ of the gradient of the local loss function f_k . The central server then aggregates those oracles and performs the update.

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

federated learning [19, Sec. 3.5]. A central idea to tackle this issue is *communication compression*, which consists in applying a lossy compression to the parameters or gradients to be transmitted. Since compression alters the message transmitted, the number of iterations required to reach a given accuracy may increase, therefore compression is of interest in situations where the communication gains are large relative to the increase of communication rounds. The design of new compression schemes (see among others [34, 2, 4, 5, 39]) and the adaptation of the learning algorithms to this setting (see e.g. [37, 1, 40, 38, 41, 25, 30, 13, 12, 23] and the references therein) are an extremely active field of research.

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

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

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

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

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

2 StVoQ algorithm

Several compression operators [39, 31, 10, 4, 8, 41, 42] have been introduced recently as bandwidth reduction for distributed learning became a major challenge. In this section, we first discuss the importance of unbiasedness of compression operators in Subsection 2.1. We then present the StVoQ compression scheme in Subsection 2.2. Finally, we compare StVoQ to competing approaches, both theoretically and empirically on a small scale example with a high compression rate.

2.1 Unbiased gradient estimate to mitigate high compression rates

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:

A1 (Unbiased Compression with relatively bounded variance). A *compression operator* Comp 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, 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 contrac-
 87 tive property along the direction of descent [37, 5, 12]:

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 compres-
 91 sion 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
 93 between the number of workers K and the compression error, a simple situation is the case in
 94 which the workers use *independent and identically distributed compression operators* $(\text{Comp}_k)_{k=1}^K$
 95 to compress the *same vector* $x \in \mathbb{R}^d$. The central node aggregates $\{\text{Comp}_k(x)\}_{k=1}^K$ into
 96 $K^{-1} \sum_{k=1}^K \text{Comp}_k(x)$. 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

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

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 vec-
 105 tor rather than each of its coordinates. A Vec-
 106 tor 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
 112 of (a) Voronoi quantization scheme before de-
 113 scribing more precisely (b) random codebooks, (c) whose distributions are invariant by unitary trans-
 114 forms. Then, (d) a method to obtain an unbiased Voronoi scheme is presented and finally (e) its
 115 asymptotic properties (as $M \rightarrow \infty$) are given.

116 **(a) Voronoi Quantization.** Voronoi quantization [26, 28], 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) [27], 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
 128 the receiver use the **same seeds** so that the same codebooks $\mathcal{C}_{M,t}$ can be reconstructed on both sides.

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 */

```

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

(c) **Unitary invariant Codewords.** Denote by $U(d) = \{U, U^*U = I\}$ the set of unitary transforms over \mathbb{R}^d . We assume in the sequel that the codeword distribution p is unitary invariant, meaning that:

A3. The distribution of the codewords p is invariant under the unitary group, i.e. for all $U \in U(d)$, and any $x \in \mathbb{R}^d$, $p(Ux) = p(x)$.

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

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

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

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

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

The proof is postponed to Appendix A.4.

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

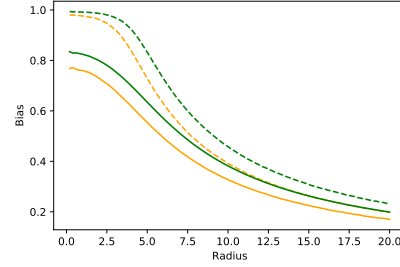


Figure 1: function r_M^p for $d = 4$ (dashed) and $d = 16$ (solid), $p = \mathcal{N}(0, I_d)$ and $M = 2^{10}$ (orange), and $M = 2^{13}$ (green).

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

A 4. (1) there exists $\epsilon > 0$ such that $\int r^{2+\epsilon} p_{\text{rad}}(r) dr < \infty$ (2) for some $\delta > 0$, $m_\delta = \inf_{r \leq \delta} p_{\text{rad}}(r) > 0$, and (3) p_{rad} is unimodal, i.e. the super level sets $\{r \in \mathbb{R}_+, p_{\text{rad}}(r) \geq t\}$, for $t \geq 0$ are convex subsets of \mathbb{R}_+ .

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

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 $x \in \mathbb{R}^d$,

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

The proof is postponed to Appendix C.1. Note that $C_d \cong_{d \rightarrow \infty} d/(2\pi e)$ hence C_d grows only linearly with the dimension d . We can now exploit this result to control the radial bias as a function of $\|x\|$. 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

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

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 $O(M^{-1/d})$. We use an a scalar quantizer SQ to transmit $1/r_M^p(\|x\|)$. Because the range of values taken by $1/r_M^p(\|x\|)$ is limited, a small number of bits P is sufficient (we typically use $P = 3$ bits). The total number of transmitted bits is $\log_2(M) + \log_2(P)$. We use a random unbiased scalar 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 elements. A scalar quantizer is said to be unbiased if $\mathbb{E}[\text{SQ}(r)] = r$ for all $r \in \mathbb{R}$. Assuming that 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 save space, we present the details of the scalar quantization (based on nonuniform random dither) methods is presented in Appendix B.1.

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

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

We stress that in this case the expectation is taken w.r.t. the input distribution q , the codebook being deterministic in (4). A *Voronoi optimal codebook* $\mathcal{C}_M^{q,*}$ is a minimizer of the distortion over 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 [14] gives the distortion of the Voronoi optimal codebook in the limit of $M \rightarrow \infty$; see Appendix C.1 for a precise 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 if $\|q\|_{d/(d+2)} < \infty$, then as $M \rightarrow \infty$, $\text{Dist}(q, \mathcal{C}_M) \cong M^{-2/d} J_d \|q\|_{d/(d+2)}$, and J_d is a universal constant J_d satisfying $J_d \cong_{d \rightarrow \infty} d/2\pi e$ (see Appendix C.2 for the exact constant).

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

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

If $\|q\|_{d/(d+2)} < \infty$, using the Hölder inequality with negative exponents (see [16, p. 191] and Appendix C.3), it holds that $C(q, p, d) \geq \|q\|_{d/(d+2)}$.

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 $\delta > 0$, and $C(q, p, d) < \infty$. Then,

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

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

The proof is postponed to Appendix C.2. In words, under general assumptions, the distortion 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 in addition q is unitarily invariant and unimodal, then a random codebook distributed according to $p_{q,d,*}$ reaches the optimal distortion bound, up to universal constants (depending only on the dimension d). Moreover, as $d \rightarrow \infty$, then $C_d J_d^{-1} \cong_{d \rightarrow \infty} 1$ and the efficiency gap vanishes. As an illustration, assume that the input distribution is standard Gaussian $q = \mathcal{N}(0, \text{I}_d)$ and set the codeword distribution to be $p_\alpha = \mathcal{N}(0, \alpha^2 \text{I}_d)$ where $\alpha^2 \in \mathbb{R}_+^*$. If $\alpha^2 d > 2$, then $C(\mathcal{N}(0, \text{I}_d), \mathcal{N}(0, \alpha^2 \text{I}_d), d) = 2\pi\alpha^2 \{\alpha^2 d / (\alpha^2 d - 2)\}^{d/2}$ and $\|\mathcal{N}(0, \text{I}_d)\|^{(2+d)/2} = (2\pi)(1 + 2/d)^{1+2/d}$. The function $\alpha \rightarrow C(\mathcal{N}(0, \text{I}_d), \mathcal{N}(0, \alpha^2 \text{I}_d), d)$ has a unique minimum at $\alpha_d^2 = 1 + 2/d$ for which $C(\mathcal{N}(0, \text{I}_d), \mathcal{N}(0, \alpha_d^2 \text{I}_d), d) = \|\mathcal{N}(0, \text{I}_d)\|^{(2+d)/2}$ showing that a random codebook sampled from $\mathcal{N}(0, \alpha_d^2 \text{I}_d)$ is optimal. It is interesting to note that the variance of the codeword distribution should be $(1 + 2/d)$ larger than the variance of the input distribution $\mathcal{N}(0, \text{I}_d)$.

2.3 Related works

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

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

Top-H/Rand H. Achieving higher compression rates is possible through *sparsification* operators, that only transmit a few coordinates. The most popular schemes are Top- H and Rand- H compressors, that respectively map the vector to either its H largest coordinates, or a random subset of cardinality H , rescaled by d/H to ensure unbiasedness. Top- H is a biased operator, and the performance of 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			StoVoQ $\log_2(M)$	DoStoVoQ $\log_2(M)$
	SGD $32d$	Sign d	QSGD $s \geq 1$ $32 + s\sqrt{d}\log_2(d)$	Top- H $32H$	Rand- H $32H$	Polytope [10] $\log_2(2d)$	HSQ-span [8] $\log_2(M)$	HSQ-greed [8] $\log_2(M)$		
Unbiased	-	-	✓	-	✓	✓	✓	-	✓	✓ (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	-	-	-	$M/\sigma_{\min}(C)$	-	-

HyperSphere Quantization (HSQ). HSQ was introduced by Dai et al. [8]. Two versions are considered: (1) a - greedy- Voronoi VQ referred to as HSQ-greed in Table 1, which is biased, and for which the theoretical guarantee provided in the paper (in their Lemma 3 and Theorem 3, which corresponds to a variant of A 2 and the subsequent convergence rate) *worsens* as M increases, making it mostly vacuous; (2) an unbiased version VQ (HSQ-span), which uses a minimum-norm decomposition of $x \in \text{Span}(\mathcal{C}_M)$ the linear subspace generated by the codewords - this version suffers from a large variance (see Table 7) and potentially an ill-conditioning. Moreover, the performance of HSQ-span does not improve with M .

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

Dual Quantization and Cross-polytope. An approach to constructing unbiased VQ is to use the dual VQ, also referred to as Delaunay Quantization (DQ); see [27]. DQ is unbiased for any $x \in \text{ConvHull}(\mathcal{C}_M)$, the convex hull of \mathcal{C}_M . DQ requires to compute the barycentric coordinates for $x \in \text{ConvHull}(\mathcal{C}_M)$, that is to solve $(\lambda_1^x, \dots, \lambda_M^x) = \arg\min_{\lambda_1, \dots, \lambda_M} \|x - \sum_{i=1}^M \lambda_i c_i\|^2$, under the constraints $\lambda_i \geq 0, \sum_{i=1}^M \lambda_i = 1$. The quantizer is obtained by drawing a codeword c_i with probability $[\lambda_1^x, \dots, \lambda_M^x]$. Computing the barycentric coordinates is in general very demanding unless \mathcal{C}_M has a very simple structure (see Appendix B for details). The Cross-Polytope method Gandikota et al. [10] is a simple instance of DQ, with a codebook $\mathcal{C}_{2d}^{\text{CP}}$ composed of the $2d$ canonical vectors $\{\pm \sqrt{d}e_i = \pm(0, \dots, 0, \sqrt{d}, 0 \dots 0), i \in [d]\}$, that relies on the inclusion $B_2(0; 1) \subset B_1(0; \sqrt{d}) = \text{ConvHull}(\mathcal{C}_{2d}^{\text{CP}})$. The barycentric decomposition can then easily be computed. Unfortunately, this method suffers from a large variance, as the quantization error $\|VQ^{\text{CP}}(x, \mathcal{C}_M) - x\|$ of any x is lower bounded by $\sqrt{d} - 1$, which means the error has the same quadratic error than the Rand-1 compressor.

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

Numerical comparisons: In Table 7, we compare the distortions achieved by the compression methods given in Table 1 for a communication budget of 16 bits for $d = 16$ and assuming that the input distribution is $q = \mathcal{N}(0, I_d)$. The compression factor is 32 (assuming 32 bits floating point per coordinate). Such a compression rate is out of reach for QSGD, that requires, even for $s = 1$ at least $\sqrt{d}\log(d) + R$ bits, where R is the number of bits to encode the norm (32 in [2]). For QSGD we have quantized the norm (using an uniform quantizer) on 3 bits and obtained an averaged distortion of 36.10 (for $K = 1$) and 1.82 for ($K = 20$) - the total number of bits is 19-. We use $H = 2$ for Top- H and Rand- H and use a scalar quantizer with 8 bits. For HSQ, we use 6 bits for the norm, using the unbiased uniform quantizer given in [8] and a Voronoi optimal codebook for the uniform distribution on the unit-sphere with $M = 2^{10}$ codewords. For StoVoQ we use a random codebook with $M = 2^{13}$ codewords; the codewords are sampled from a $\mathcal{N}(0, (1 + 2/d)I_d)$, and 3 bits are allocated for the scalar quantization of $1/r_M^p$ (the inverse of the radial bias). Finally, we average the result of 2 independent compressions for Polytope (following the replication technique described in [10]). We use $n = 10^4$ vectors, and report in Table 7 the distortion and sample variance. For StoVoQ 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)

3 DoStoVoQ algorithm

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

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

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

Bias and variance of the compressed gradient with K workers.

Consider the two filtrations $(\mathcal{F}_t)_{t \geq 0}$ and $(\mathcal{G}_t)_{t \geq 0}$ defined recursively as follows $\mathcal{F}_0 = \sigma(\emptyset)$ and for $t \geq 0$, $\mathcal{G}_{t+1} = \mathcal{F}_t \vee \sigma(\{g_{k,t+1}, k \in [K]\})$ and $\mathcal{F}_{t+1} = \mathcal{G}_{t+1} \vee \sigma(\{\hat{g}_{k,t+1}, k \in [K]\})$. With these notations, for any $t \geq 0$, θ_t is \mathcal{F}_t -measurable.

Theorem 4. *At any iteration $t + 1$ in DoStoVoQ, the K compressed stochastic gradients $(\hat{g}_{k,t+1})_{k \in [K]}$ are (i) independent conditionally to \mathcal{G}_{t+1} (ii) conditionally unbiased, i.e., for all $k \in [K]$, we have $\mathbb{E}[\hat{g}_{k,t+1} | \mathcal{G}_{t+1}] = g_{k,t+1}$, (iii) satisfy the relatively bounded error condition of A 1, i.e. there exists a constant ω_M such that, for all $k \in [K]$: $\mathbb{E}[\|\hat{g}_{k,t+1} - g_{k,t+1}\|^2 | \mathcal{G}_{t+1}] \leq \omega_M \|g_{k,t+1}\|^2$.*

Moreover, ω_M decreases with the number of codewords M and the P , as $\omega_M = O(M^{-2/d}) + O(2^{-P})$ [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  $\theta_{t-1}$  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  $\theta_{t-1}$ ;
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
```


The first statement stems from the fact that each bucket is quantized using StoVoQ which is unbiased. The second statement is more challenging; proof is postponed to Appendix A.5. We stress that this result differs from Theorem 2, which corresponds to the distortion of a source with distribution q .

Convergence results. Theorem 4 proves that our compression method satisfies the assumptions needed to obtain fast convergence rate, for DoStoVoQ-SGD, and for its variants DoStoVoQ-(VR)-DIANA. Consider a Smooth and Strongly Convex (SSC) function $F = \sum_{k=1}^K f_k$, with condition number $\kappa > 1$. We measure the complexity of the algorithm by the number of iterations t required 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 [17], which provides a complexity of $O_{\kappa \rightarrow \infty}(\kappa(1 + \omega_M/K) \log(\epsilon^{-1}))$ [17, Corollary 2], applies to DoStoVoQ-VR-DIANA.

Convergence rates for DoStoVoQ-DIANA (without VR), and on non-convex optimization problems can be obtained from Horváth et al. [17, Corollary 1,3,4]. As in the strongly-convex case, complexities increase by a factor depending on $(1 + \omega_M/K)$ w.r.t. uncompressed algorithm. Intuitively, *the impact on the optimization complexity of a high compression is mitigated by the number of workers*, which supports the use of independent and unbiased compressors when the number of workers is large and high compression factors are required.

Indeed, these complexities can be compared to: (1) the one of *uncompressed* variance reduced distributed methods [9] that achieve a complexity of $O_{\kappa \rightarrow \infty}(\kappa \log(\epsilon^{-1}))$ (in the SSC case); (2) the complexity for biased compression operators satisfying A 2, Beznosikov et al. [5, Theorem 13] that obtain $O_{\kappa \rightarrow \infty}(\kappa(1 + \delta) \log(\epsilon^{-1}))$ for compressed GD (independently of the number of workers); (3) the complexities of compressed SGD methods with *error feedback* in [12]², that also have no dependency on the number of workers. **Overall, the unbiased character is crucial to mitigate the variance increase resulting from high compression rates.**

4 Numerical experiments

4.1 Least Squares Regression (LSR)

We consider a least-squares problem with $n = 2^{14}$ samples, a bucket size $d = 16$, $D = 2^9$, and $K = 32$ workers; each worker has access to a subset $m = 2^{11}$ samples (picked with replacement) to introduce a dependency in the data used by the workers. For $i \in [n]$, we assume $X_i \sim \mathcal{N}(0, I_D)$ and $Y_i \sim \mathcal{N}(X_i^\top \omega_*, 1)$ where $\omega_* \in \mathbb{R}^D$. We solve $\inf_{\omega \in \mathbb{R}^D} \sum_{i=1}^n \|Y_i - X_i^\top \omega\|^2$ via a gradient descent with step size $1/\alpha L$ where α is fine-tuned for each quantization method and $L \approx 2n$ is the smoothness constant. We use DoStoVoQ with $M = 2^{13}$ codewords sampled from $\mathcal{N}(0, (1+2/d)I_d)$ for DoStoVoQ and $M = 2^{10}$ on the unit Sphere for HSQ s.t. the number of bits transmitted at each round by the worker is set to 16 (see Table 7).

Figure 2 reports the excess-log of the train loss over $T = 10$ iterations, for a standard GD. DoStoVoQ outperforms HSQ-greed: indeed the linear convergence rate of distributed GD is faster for an unbiased compressor than for the biased approach.

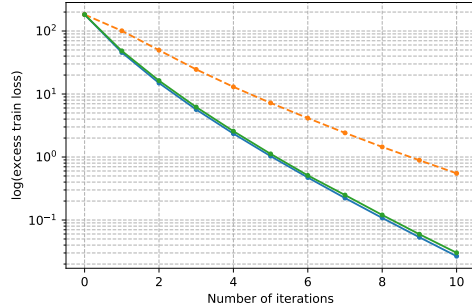


Figure 2: Comparison between GD (blue), HSQ-greed (orange) and DoStoVoQ (green), on a LSR problem in dimension $D = 2^9$.

4.2 Applications to Deep Neural Networks training

Setting. We now describe our experimental framework for training two standard models of Deep Neural Networks: a VGG-16 [35] and a ResNet-18 [15]. We follow the standard procedure of training those models both on CIFAR-10 and ImageNet; the hyper-parameters are fine-tuned to optimize the accuracy *without quantization*. We do not compress the affine constant part of the affine convolutional layers and batch normalization layers. We apply independent DoStoVoQ on

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

361 batches of 32 buckets of size $d = 16$ (i.e. we transmit a high-resolution norm for $D = 32 \cdot 16 = 512$
 362 coefficients).

363 **CIFAR-10.** We use the implementation of HSQ [8]: the batch size is 256 for CIFAR-10, the
 364 total number of epochs is 100, the initial learning rate is 0.1, which is divided by 10 and 50 at
 365 epochs 51 and 71. We report the accuracy of DoStoVoQ, QSGD, and HSQ-greed in table 4. By
 366 design, the compression factor of Q-SGD for $d = 16$ is 13, which is significantly less than HSQ
 367 or DoStoVoQ. Both HSQ and DoStoVoQ perform similarly and the accuracy gap between the two
 368 methods are under the sample variance (computed over 5 seed and about 0.2). In Table 4 we report
 369 the distortion of a random subset of gradients $\mathcal{G} = \{g_t, t \in [|\mathcal{G}|]\}$ (with $|\mathcal{G}| = 10^2$, $d = 16$, $D =$
 370 $2^5 \times d$) obtained from a given layer of a VGG on CIFAR-10, i.e.: $|\mathcal{G}|^{-1} \sum_{g_t \in \mathcal{G}} \|K^{-1} \sum_{k=1}^K (g_{k,t} -$
 371 $\hat{g}_{k,t})\|^2$, where $(\hat{g}_{k,t})_{k \in [K]}$ correspond to k independent workers compressing their own gradient
 372 $g_{k,t}$. The choice of the layer does not affect significantly the results. Even with the actual gradient
 373 distribution, DoStoVoQ outperforms for a given compression factor each unbiased method. This is
 374 on pair with the observation that the gradients of a Deep Neural Network are approximately Gaussian
 375 distributed [3, 41, 4]. Additional experiments can be found in the Appendix.

376 **ImageNet.** For ImageNet, we use different bucket sizes, the standard batch size of 256, and only
 377 $K = 1$ worker for energy savings (recall Imagenet training last about 1 day for a single worker on
 378 academic hardware). An initial learning rate of 0.1 is divided by 10 at epoch 30 and 60, while the
 379 model is trained for 90 epochs. A ResNet here obtains 69.9%, and with a compression factor of 8,
 380 the performance drops by 2.5%. Using $d = 16$, we reach a compression factor of 38, while the Top-
 381 1 accuracy drops by only 4.8%: this is a substantially higher compression rate than the concurrent
 382 work QSGD on the ImageNet dataset.

383 **Computational impact.** In the case of deep Neural Networks, our training procedure requires
 384 neither a substantial modifications of standard pipelines, nor a modification of the hyper-parameters
 385 which allows to save computational resources. Green Algorithm ([22]) shows that this work gen-
 386 erated around 15kg of CO2, and require 400 kWh. A typical experiment lasted few hours on CIFAR-
 387 10 and about 3 days on ImageNet, which is in the standard range for this type of prototypical codes.
 388 This work could have future impact on FL, to reduce their electrical consumption.

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

References

- [1] Naman Agarwal, Ananda Theertha Suresh, Felix Xinnan X Yu, Sanjiv Kumar, and Brendan McMahan. cpSGD: Communication-efficient and differentially-private distributed SGD. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 7564–7575. Curran Associates, Inc., 2018.
- [2] Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. QSGD: Communication-Efficient SGD via Gradient Quantization and Encoding. *Advances in Neural Information Processing Systems*, 30:1709–1720, 2017.
- [3] Ron Banner, Itay Hubara, Elad Hoffer, and Daniel Soudry. Scalable methods for 8-bit training of neural networks. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 5151–5159, 2018.
- [4] Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. signsgd: Compressed optimisation for non-convex problems. In *International Conference on Machine Learning*, pages 560–569. PMLR, 2018.
- [5] Aleksandr Beznosikov, Samuel Horváth, Peter Richtárik, and Mher Safaryan. On Biased Compression for Distributed Learning. *arXiv:2002.12410 [cs, math, stat]*, February 2020. arXiv: 2002.12410.
- [6] Léon Bottou. On-line learning and stochastic approximations. 1999. doi: 10.1017/CBO9780511569920.003.
- [7] Léon Bottou. Large-Scale Machine Learning with Stochastic Gradient Descent. In Yves Lechevallier and Gilbert Saporta, editors, *Proceedings of COMPSTAT’2010*, pages 177–186, Heidelberg, 2010. Physica-Verlag HD. ISBN 978-3-7908-2604-3. doi: 10.1007/978-3-7908-2604-3_16.
- [8] Xinyan Dai, Xiao Yan, Kaiwen Zhou, Han Yang, Kelvin KW Ng, James Cheng, and Yu Fan. Hyper-sphere quantization: Communication-efficient sgd for federated learning. *arXiv preprint arXiv:1911.04655*, 2019.
- [9] Aaron Defazio, Francis R Bach, and Simon Lacoste-Julien. Saga: A fast incremental gradient method with support for non-strongly convex composite objectives. In *NIPS*, 2014.
- [10] Venkata Gandikota, Daniel Kane, Raj Kumar Maity, and Arya Mazumdar. vqsgd: Vector quantized stochastic gradient descent. In *International Conference on Artificial Intelligence and Statistics*, pages 2197–2205. PMLR, 2021.
- [11] Richard Gardner. The brunn-minkowski inequality. *Bulletin of the American Mathematical Society*, 39(3):355–405, 2002.
- [12] Eduard Gorbunov, Dmitry Kovalev, Dmitry Makarenko, and Peter Richtárik. Linearly Converging Error Compensated SGD. *arXiv:2010.12292 [cs, math]*, October 2020. arXiv: 2010.12292.
- [13] Eduard Gorbunov, Konstantin Burlachenko, Zhize Li, and Peter Richtárik. Marina: Faster non-convex distributed learning with compression. *arXiv preprint arXiv:2102.07845*, 2021.
- [14] Siegfried Graf and Harald Luschgy. *Foundations of quantization for probability distributions*. Springer, 2007.
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [16] Edwin Hewitt and Karl Stromberg. *Real and abstract analysis: a modern treatment of the theory of functions of a real variable*. Springer-Verlag, 2013.

- [17] Samuel Horváth, Dmitry Kovalev, Konstantin Mishchenko, Sebastian Stich, and Peter Richtárik. Stochastic Distributed Learning with Gradient Quantization and Variance Reduction. *arXiv:1904.05115 [math]*, April 2019. arXiv: 1904.05115.
- [18] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 1*, NIPS’13, pages 315–323, Lake Tahoe, Nevada, December 2013. Curran Associates Inc.
- [19] Peter Kairouz, H. Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Keith Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, Rafael G. L. D’Oliveira, Salim El Rouayheb, David Evans, Josh Gardner, Zachary Garrett, Adrià Gascón, Badih Ghazi, Phillip B. Gibbons, Marco Gruteser, Zaid Harchaoui, Chaoyang He, Lie He, Zhouyuan Huo, Ben Hutchinson, Justin Hsu, Martin Jaggi, Tara Javidi, Gauri Joshi, Mikhail Khodak, Jakub Konečný, Aleksandra Korolova, Farinaz Koushanfar, Sanmi Koyejo, Tancrède Lepoint, Yang Liu, Prateek Mittal, Mehryar Mohri, Richard Nock, Ayfer Özgür, Rasmus Pagh, Mariana Raykova, Hang Qi, Daniel Ramage, Ramesh Raskar, Dawn Song, Weikang Song, Sebastian U. Stich, Ziteng Sun, Ananda Theertha Suresh, Florian Tramèr, Praneeth Vepakomma, Jianyu Wang, Li Xiong, Zheng Xu, Qiang Yang, Felix X. Yu, Han Yu, and Sen Zhao. Advances and Open Problems in Federated Learning. *arXiv:1912.04977 [cs, stat]*, December 2019. arXiv: 1912.04977.
- [20] Jakub Konečný, H. Brendan McMahan, Daniel Ramage, and Peter Richtárik. Federated Optimization: Distributed Machine Learning for On-Device Intelligence. *arXiv:1610.02527 [cs]*, October 2016. arXiv: 1610.02527.
- [21] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [22] Loïc Lannelongue, Jason Grealey, and Michael Inouye. Green Algorithms: Quantifying the carbon emissions of computation. *arXiv:2007.07610 [cs]*, October 2020. arXiv: 2007.07610.
- [23] Zhize Li, Dmitry Kovalev, Xun Qian, and Peter Richtarik. Acceleration for Compressed Gradient Descent in Distributed and Federated Optimization. In *International Conference on Machine Learning*, pages 5895–5904. PMLR, November 2020. ISSN: 2640-3498.
- [24] Tao Lin, Sebastian U. Stich, Kumar Kshitij Patel, and Martin Jaggi. Don’t Use Large Mini-Batches, Use Local SGD. *arXiv e-prints*, art. arXiv:1808.07217, August 2018.
- [25] Konstantin Mishchenko, Eduard Gorbunov, Martin Takáč, and Peter Richtárik. Distributed Learning with Compressed Gradient Differences. *arXiv:1901.09269 [cs, math, stat]*, June 2019. arXiv: 1901.09269.
- [26] Gilles Pagès and Jacques Printems. Optimal quadratic quantization for numerics: the gaussian case. *Monte Carlo methods and applications*, 9(2):135–165, 2003.
- [27] Gilles Pagès and Benedikt Wilbertz. Sharp rate for the dual quantization problem. In *Séminaire de Probabilités XLIX*, volume 2215 of *Lecture Notes in Math.*, pages 405–454. Springer, Cham, 2018.
- [28] Gilles Pagès and Benedikt Wilbertz. Sharp rate for the dual quantization problem. In *Séminaire de Probabilités XLIX*, pages 405–454. Springer, 2018.
- [29] Abhishek Panigrahi, Raghav Somani, Navin Goyal, and Praneeth Netrapalli. Non-gaussianity of stochastic gradient noise. *arXiv preprint arXiv:1910.09626*, 2019.
- [30] Constantin Philippenko and Aymeric Dieuleveut. Artemis: tight convergence guarantees for bidirectional compression in Federated Learning. *arXiv:2006.14591 [cs, stat]*, November 2020. arXiv: 2006.14591.
- [31] Ali Ramezani-Kebrya, Fartash Faghri, and Daniel M Roy. Nuqsgd: Improved communication efficiency for data-parallel sgd via nonuniform quantization. *arXiv preprint arXiv:1908.06077*, 2019.

- [32] Herbert Robbins and Sutton Monro. A Stochastic Approximation Method. *Annals of Mathematical Statistics*, 22(3):400–407, September 1951. ISSN 0003-4851, 2168-8990. doi: 10.1214/aoms/1177729586. Number: 3 Publisher: Institute of Mathematical Statistics.
- [33] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- [34] F. Seide, H. Fu, Jasha Droppo, G. Li, and D. Yu. 1-bit stochastic gradient descent and its application to data-parallel distributed training of speech DNNs. pages 1058–1062, January 2014.
- [35] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [36] Sebastian U. Stich. Local SGD Converges Fast and Communicates Little. *arXiv:1805.09767 [cs, math]*, May 2019. arXiv: 1805.09767.
- [37] Sebastian U. Stich and Sai Praneeth Karimireddy. The Error-Feedback Framework: Better Rates for SGD with Delayed Gradients and Compressed Communication. *arXiv:1909.05350 [cs, math, stat]*, September 2019. arXiv: 1909.05350.
- [38] Sebastian U Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified SGD with Memory. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems 31*, pages 4447–4458. Curran Associates, Inc., 2018.
- [39] Thijs Vogels, Sai Praneeth Karimireddy, and Martin Jaggi. Powersgd: Practical low-rank gradient compression for distributed optimization. *Advances in Neural Information Processing Systems*, 32:14259–14268, 2019.
- [40] Jianqiao Wangni, Jialei Wang, Ji Liu, and Tong Zhang. Gradient Sparsification for Communication-Efficient Distributed Optimization. *Advances in Neural Information Processing Systems*, 31:1299–1309, 2018.
- [41] An Xu, Zhouyuan Huo, and Heng Huang. Optimal gradient quantization condition for communication-efficient distributed training. *arXiv preprint arXiv:2002.11082*, 2020.
- [42] Hantian Zhang, Jerry Li, Kaan Kara, Dan Alistarh, Ji Liu, and Ce Zhang. The zipml framework for training models with end-to-end low precision: The cans, the cannots, and a little bit of deep learning. *arXiv preprint arXiv:1611.05402*, 2016.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#) See Section 2 for quantization and Section 4 for associated experiments.
 - (b) Did you describe the limitations of your work? [\[Yes\]](#) See broader impact and Appendix.
 - (c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#) Detailed experiments carbon footprint can be find in Section 4.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#)
 - (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#) Also see Appendix in Supplemental Material.
3. If you ran experiments...

- 536 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
537 mental results (either in the supplemental material or as a URL)? [Yes] Code available
538 in Supplementary Material.
- 539 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
540 were chosen)? [Yes] See Section 4.
- 541 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
542 ments multiple times)? [Yes] In particular Table 7 presents standard deviations, and
543 variances of NN model accuracies from Section 4 can be found in Appendix.
- 544 (d) Did you include the total amount of compute and the type of resources used (e.g.,
545 type of GPUs, internal cluster, or cloud provider)? [Yes] See Section 4 for further
546 references.
- 547 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 548 (a) If your work uses existing assets, did you cite the creators? [Yes] As mentioned in
549 Section 4, code is partly inspired from [8].
- 550 (b) Did you mention the license of the assets? [Yes] Only open source and/or Academic
551 assets are used.
- 552 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
553 Radial biases already computed available in Supplementary Material.
- 554 (d) Did you discuss whether and how consent was obtained from people whose data
555 you’re using/curating? [N/A] Use of publicly available data (CIFAR10 [21] and Im-
556 genet [33]).
- 557 (e) Did you discuss whether the data you are using/curating contains personally identifi-
558 able information or offensive content? [N/A]
- 559 5. If you used crowdsourcing or conducted research with human subjects...
- 560 (a) Did you include the full text of instructions given to participants and screenshots, if
561 applicable? [N/A]
- 562 (b) Did you describe any potential participant risks, with links to Institutional Review
563 Board (IRB) approvals, if applicable? [N/A]
- 564 (c) Did you include the estimated hourly wage paid to participants and the total amount
565 spent on participant compensation? [N/A]

566 Contents

567	A Proofs	14
568	A.1 Classical compressors mentioned in the main text	14
569	A.2 Notations	14
570	A.3 Proof of Lemma 1	14
571	A.4 Proof of Theorem 1	14
572	A.5 Proof of Theorem 4	15
573	B Scalar and vector Quantization	19
574	B.1 Unbiased random scalar quantization	19
575	B.2 Dual Vector Quantization	20
576	B.3 HSQ methods - see Dai et al. [8]	20
577	B.4 Alignment under A 3, for StoVoQ without debiasing function (new)	23
578	C Unitarily invariant random codebooks	23
579	C.1 Proof of Theorem 2	23
580	C.2 Proof of Theorem 3	25
581	C.3 An elementary lower-bound	25
582	C.4 Asymptotic distortion of a random quantizer on the unit sphere S_{d-1}	25
583	D Algorithmic extensions	26
584	D.1 Spherical codebooks	26
585	D.2 Extension to DoStoVoQ-DIANA and DoStoVoQ-VR-DIANA	26

586	E Additional experiments	28
587	E.1 Distortion for Gaussian input	28
588	E.2 Distortion for neural networks gradients	29

589 A Proofs

590 A.1 Classical compressors mentioned in the main text

591 For completeness, we here recall the formal definitions of the scalar compression operators men-
592 tioned in the main text. For $i \in [d]$, denote by e_i the i -th canonical vector. Let $H \in [d]$.

593 **Definition 1 (Sign).** For any $x \in \mathbb{R}^d$, $\text{Sign}(x) := \sum_{i \in [d]} \text{sign}(x_i) e_i$.

594 **Definition 2 (Top-H).** For any $x \in \mathbb{R}^d$, $\text{Top-H}(x) := \sum_{i \in T_H} x_i e_i$, where T_H is the set composed
595 of the indices of the H largest (in absolute value) coordinates of x .

596 **Definition 3 (Rand-H).** For any $x \in \mathbb{R}^d$, $\text{Rand-H}(x) := \frac{d}{H} \sum_{i \in R_H} x_i e_i$, where R_H is the set
597 composed of H random indices picket uniformly without replacement.

598 **Definition 4 (s -quantization operator).** Let $s \geq 1$ and $p \geq 1$. Given $x \in \mathbb{R}^d$, the s -quantization
599 operator \mathcal{C}_s is defined by:

$$\mathcal{C}_s(x) := \|x\|_p \times \sum_{i=1}^d \text{sign}(x_i) \{s^{-1} \lfloor s|x_j|/\|x\|_p \rfloor + \mathbb{1}(\{U_i \leq s|x_j|/\|x\|_p - \lfloor s|x_j|/\|x\|_p \rfloor\})\} e_i.$$

600 where $\{U_i\}_{i=1}^d$ are d -independent uniform random variables on $[0, 1]$.

601 The s -quantization scheme verifies A 1 with $\omega_s = \min(d/s^2, \sqrt{d}/s)$. Proof can be found in Alistarh
602 et al. [2, see Appendix A.1].

603 A.2 Notations

604 For $u, v \in \mathbb{R}^d$, $\langle u, v \rangle = u^\top v$ denotes the standard scalar product. For $p \geq 1$ and $x \in \mathbb{R}^d$, $\|x\|_p =$
605 $\left\{ \sum_{i=1}^d |x_i|^p \right\}^{1/p}$. When $p = 2$, we sometimes drop the subscript, i.e. we write $\|x\|$ as a shorthand
606 notation of $\|x\|_2$.

607 A function $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}$ is said to be a *radial function* if and only if φ is invariant under unitary
608 transforms, i.e. for all $x \in \mathbb{R}^d$ and $U \in \text{U}(d)$, $\varphi(Ux) = \varphi(x)$.

609 We denote for $t > 0$ by $\Gamma(t) = \int_0^{+\infty} u^{t-1} e^{-u} du$ the Gamma function. Leb d the Lebesgue measure
610 on \mathbb{R}^d . $B(x; r)$ is the (Euclidean) ball centered at $x \in \mathbb{R}^d$ with radius $r > 0$. We denote by
611 $S_{d-1} = \{x \in \mathbb{R}^d, \|x\| = 1\}$ the unit-sphere and σ_{d-1} the uniform distribution ofn S_{d-1} .

612 A.3 Proof of Lemma 1

613 Note that, for any $U \in \text{U}(d)$ and $x \in \mathbb{R}^d$,

$$\text{VQ}(Ux, \mathcal{C}_M) = \arg\min_{c \in \mathcal{C}_M} \|Ux - c\| = \arg\min_{c \in \mathcal{C}_M} \|x - U^\top c\| = U \text{VQ}(x, U^\top \mathcal{C}_M), \quad (7)$$

614 where $U^\top \mathcal{C}_M = \{U^\top C_1, \dots, U^\top C_n\}$. Using (7) and A 3, we get

$$\mathbb{E}_{\mathcal{C}_M \sim p}[g(\text{VQ}(Ux, \mathcal{C}_M))] = \mathbb{E}_{\mathcal{C}_M \sim p}[g(U \text{VQ}(x, U^\top \mathcal{C}_M))] = \mathbb{E}_{\mathcal{C}_M \sim p}[g(U \text{VQ}(x, \mathcal{C}_M))].$$

615 A.4 Proof of Theorem 1

616 We preface the proof of the Theorem by stating and proving two elementary lemmas.

617 **Lemma 2.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be a function such that $f(Ux) = Uf(x)$ for any $x \in S_{d-1}$ and
618 $U \in \text{U}(d)$. Then, there exists $r \in \mathbb{R}$ such that $f(x) = rx$ for all $x \in \mathbb{R}^d$.

619 *Proof.* For all $x \in S_{d-1}$, define $g(x) = f(x) - \langle f(x), x \rangle x$. It is easily checked that for all $x \in \mathbb{R}^d$
620 and $U \in U(d)$, $g(Ux) = Ug(x)$. Let U_x be the reflection symmetry with axis $\mathbb{R}x$: $U_x x = x$ and
621 for any vector $y \in \mathbb{R}^d$ orthogonal to x , $U_x y = -y$. Since $g(x) = g(U_x x) = U_x g(x) = -g(x)$, we
622 get that $g(x) = 0$ for all $x \in S_{d-1}$. Finally, denote by $U_{x \rightarrow e_1}$ (where e_1 is the first canonical vector)
623 any unitary transform satisfying $U_{x \rightarrow e_1} x = e_1$. We get

$$\langle f(x), x \rangle = \langle U_{x \rightarrow e_1}^\top f(U_{x \rightarrow e_1} x), x \rangle = \langle f(U_{x \rightarrow e_1} x), U_{x \rightarrow e_1} x \rangle = \langle f(e_1), e_1 \rangle = r,$$

624 which concludes the proof. \square

625 **Lemma 3.** Let $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ be a function such that $f(Ux) = Uf(x)$ for any $x \in \mathbb{R}^d$ and
626 $U \in U(d)$. Then, there exists a function $r : \mathbb{R} \rightarrow \mathbb{R}$ such that $f(x) = r(\|x\|)x$. Moreover,
627 $r(x) = \|f(\|x\|e_1)\|/\|x\|$.

628 *Proof.* Let $\lambda > 0$, and define for $x \in S_{d-1}$, $f_\lambda(x) = f(\lambda x)$. Lemma 2 shows that there exists
629 $\rho(\lambda) \in \mathbb{R}$ such that, for all $x \in S_{d-1}$, $f_\lambda(x) = f(\lambda x) = \rho(\lambda)x$. Hence for $x \in \mathbb{R}^d$, $f(x) =$
630 $f_{\|x\|}(x/\|x\|) = \rho(\|x\|)x/\|x\|$. Hence $|\rho(\|x\|)| = \|f(x)\| = \|f(\|x\|U_{x/\|x\| \rightarrow e_1} x/\|x\|)\| =$
631 $\|f(\|x\|e_1)\|$. The proof follows. \square

632 *Proof of Theorem 1.* The proof follows from Lemmas 1 and 3. \square

633 A.5 Proof of Theorem 4

634 We preface the proof by several technical lemmas. These lemmas establish important properties of
635 random vector quantization that are of interest beyond the proof of the theorem.

636 **Lemma 4.** Let $c_1, c_2, x \in \mathbb{R}^d$ and $\lambda \in (0, 1)$. If $\|x - c_1\| \leq \|x - c_2\|$ and $\|\lambda x - c_2\| \leq \|\lambda x - c_1\|$,
637 then $\|c_2\| \leq \|c_1\|$.

638 *Proof.* Indeed, we have both $\|x\|^2 - 2\langle x, c_1 \rangle + \|c_1\|^2 \leq \|x\|^2 - 2\langle x, c_2 \rangle + \|c_2\|^2$ and $\lambda^2\|x\|^2 -$
639 $2\lambda\langle x, c_2 \rangle + \|c_2\|^2 \leq \lambda^2\|x\|^2 - 2\lambda\langle x, c_1 \rangle + \|c_1\|^2$. Thus $-2\langle x, c_1 \rangle \leq -2\langle x, c_2 \rangle + \|c_2\|^2 - \|c_1\|^2$
640 and $-2\langle x, c_2 \rangle \leq -2\langle x, c_1 \rangle + \lambda^{-1}\|c_1\|^2 - \lambda^{-1}\|c_2\|^2$. Combining both inequalities, we get $(\lambda^{-1} -$
641 $1)\|c_2\|^2 \leq (\lambda^{-1} - 1)\|c_1\|^2$ and as $\lambda^{-1} - 1 > 0$, we conclude $\|c_2\|^2 \leq \|c_1\|^2$. \square

642 We now make an additional assumption on the codeword distribution p .

643 **A5.** The distribution p is radially homogeneous, i.e. p is unitarily invariant and for any $\beta \in (0, 1]$
644 and $x \in \mathbb{R}^d$:

$$\|\mathbb{E}_{\mathcal{C}_M \sim p}[\text{VQ}(x, \beta \mathcal{C}_M)]\| \leq \|\mathbb{E}_{\mathcal{C}_M \sim p}[\text{VQ}(x, \mathcal{C}_M)]\|$$

645 In words, this means that contracting all codewords by a factor $\beta \in (0, 1]$ reduces the norm of the
646 expectation of the nearest neighbor of any x . This condition is slightly more restrictive than A 3.
647 It is satisfied by the standard Gaussian distribution. Under this assumption, we have the following
648 Lemma.

649 **Lemma 5.** Assume A 3-A 5 and consider the function $r_M^p : \mathbb{R}_+ \rightarrow \mathbb{R}$ defined in Theorem 1. For any
650 $M \in \mathbb{N}^*$ the function $\rho \rightarrow r_M^p(\rho)$ is non-increasing on \mathbb{R}_+ . If in addition A 4 is satisfied, then for
651 any $\rho \in \mathbb{R}_+^*$, $r_M^p(\rho) \leq 1$.

652 *Proof.* Let $x \in \mathbb{R}^d$ and $\lambda > 1$. By definition, we have:

$$r_M^p(\lambda\|x\|)\lambda x = \mathbb{E}_{\mathcal{C}_M \sim p}[\text{VQ}(\lambda x, \mathcal{C}_M)] = \lambda \mathbb{E}_{\mathcal{C}_M \sim p}[\text{VQ}(x, \lambda^{-1} \mathcal{C}_M)], \quad (8)$$

653 where we have used that a.s., $\text{VQ}(\lambda x, \mathcal{C}_M) = \lambda \text{VQ}(x, \lambda^{-1} \mathcal{C}_M)$. On the other hand:

$$\mathbb{E}_{\mathcal{C}_M \sim p}[\text{VQ}(x, \mathcal{C}_M)] = r_M^p(\|x\|)x. \quad (9)$$

654 Combining both equations under A 2, we get:

$$\begin{aligned} \|r_M^p(\lambda\|x\|)\lambda x\|^2 &\stackrel{\text{eq. (8)}}{=} \|\lambda \mathbb{E}_{\mathcal{C}_M \sim p}[\text{VQ}(x, \lambda^{-1} \mathcal{C}_M)]\|^2 \\ &\stackrel{\text{A 5}}{\leq} \|\lambda \mathbb{E}_{\mathcal{C}_M \sim p}[\text{VQ}(x, \mathcal{C}_M)]\|^2 \\ &\stackrel{\text{eq. (9)}}{\leq} \|\lambda r_M^p(\|x\|)x\|^2. \end{aligned}$$

Overall, we obtain that $(r_M^p(\lambda\|x\|))^2 \leq (r_M^p(\|x\|))^2$, thus that r_M^p is non-increasing.

We now consider the second statement. Note first that

$$\mathbb{E}_{\mathcal{C}_M \sim p} \left[\max_{i \in [M]} \|C_i\| \right] \leq \mathbb{E}_{\mathcal{C}_M \sim p} \left[\sum_{i=1}^M \|C_i\|^2 \right] = M \mathbb{E}_{C_1 \sim p} [\|C_1\|^2]. \quad (10)$$

Since $\|\text{VQ}(x, \mathcal{C}_M)\| \leq \max_{i \in [M]} \|C_i\|$, for all $x \in \mathbb{R}^d$ it holds that

$$\|\mathbb{E}_{\mathcal{C}_M \sim p} [\text{VQ}(x, \mathcal{C}_M)]\| \leq M^{1/2} (\mathbb{E}_{C_1 \sim p} [\|C_1\|^2])^{1/2}.$$

Hence, for all $x \in \mathbb{R}^d$ such that $\|x\| \geq M^{1/2} (\mathbb{E}_{C_1 \sim p} [\|C_1\|^2])^{1/2}$, $\|\mathbb{E}_{\mathcal{C}_M \sim p} [\text{VQ}(x, \mathcal{C}_M)]\| \leq \|x\|$. For all $\lambda \in (0, 1)$, using A 5, we get

$$r_M^p(\lambda\|x\|)\lambda\|x\| = \|\mathbb{E}_{\mathcal{C}_M \sim p} [\text{VQ}(\lambda x, \mathcal{C}_M)]\| = \lambda \|\mathbb{E}_{\mathcal{C}_M \sim p} [\text{VQ}(x, \lambda^{-1} \mathcal{C}_M)]\| \leq \lambda\|x\|,$$

which concludes the proof. \square

Lemma 6. Assume A 3-A 4. Then, for any $M \in \mathbb{N}^*$, $\mathbb{E}_{\mathcal{C}_M \sim p} [\|\text{VQ}(x, \mathcal{C}_M) - x\|^2]$ is (a radial function) which is non-decreasing, i.e. for any $x \in \mathbb{R}^d$ and $\lambda \in [0, 1]$, $\mathbb{E}_{\mathcal{C}_M \sim p} [\|\text{VQ}(\lambda x, \mathcal{C}_M) - \lambda x\|^2] \leq \mathbb{E}_{\mathcal{C}_M \sim p} [\|\text{VQ}(x, \mathcal{C}_M) - x\|^2]$.

Proof. By Lemma 1, for any $U \in \text{U}(d)$ and $x \in \mathbb{R}^d$, we get

$$\mathbb{E}_{\mathcal{C}_M \sim p} [\|\text{VQ}(Ux, \mathcal{C}_M) - Ux\|^2] = \mathbb{E}_{\mathcal{C}_M \sim p} [\|U \text{VQ}(x, \mathcal{C}_M) - Ux\|^2] = \mathbb{E}_{\mathcal{C}_M \sim p} [\|\text{VQ}(x, \mathcal{C}_M) - x\|^2]$$

showing that this function is radial. We write, for any $x \in \mathbb{R}^d$:

$$\begin{aligned} \mathbb{E}_{\mathcal{C}_M \sim p} [\|\text{VQ}(x, \mathcal{C}_M) - x\|^2] &= 2 \int_{t=0}^{\infty} t \mathbb{P}_{\mathcal{C}_M \sim p} (\|\text{VQ}(x, \mathcal{C}_M) - x\|^2 > t) dt \\ &= 2 \int_{t=0}^{\infty} t \mathbb{P}_{\mathcal{C}_M \sim p} \left(\min_{i \in [M]} \|C_i - x\|^2 > t \right) dt \\ &= 2 \int_{t=0}^{\infty} t (1 - \mathbb{P}_{C_1 \sim p} (\|C_1 - x\|^2 \leq t))^M dt \\ &= 2 \int_{t=0}^{\infty} t \left(1 - \mathbb{P}_{C_1 \sim p} (B_2(x; \sqrt{t})) \right)^M dt. \end{aligned}$$

By Anderson's theorem [11], we have that $\mathbb{P}_{C_1 \sim p} (B_2(x; \sqrt{t}))$ is (a radial function) which is non-increasing, i.e. for any $x \in \mathbb{R}^d$ and $\lambda \in [0, 1]$, $\mathbb{P}_{C_1 \sim p} (B_2(\lambda x; \sqrt{t})) \geq \mathbb{P}_{C_1 \sim p} (B_2(x; \sqrt{t}))$.

Consequently, the quadratic error $\mathbb{E}_{\mathcal{C}_M \sim p} [\|\text{VQ}(x, \mathcal{C}_M) - x\|^2]$ is non-decreasing radial function. \square

Next, we provide a control on the second order moment of $\text{VQ}(x, \mathcal{C}_M)$.

Lemma 7. Assume A 3-A 4. Then, for any $M \in \mathbb{N}^*$, the function $x \mapsto \mathbb{E}_{\mathcal{C}_M \sim p} [\|\text{VQ}(x, \mathcal{C}_M)\|^2] =: M_{d,p,M}(\|x\|)$ is radial and $r \mapsto M_{d,p,M}(r)$ is non-decreasing. Moreover, for any $M_0 \geq 1$ there exists a constant $C_{M_0,R,d,p}$ such that for all $M \geq M_0$, $M_{d,p,M}(R) \leq C_{M_0,R,d,p}$.

Proof. The fact that $x \mapsto \mathbb{E}_{\mathcal{C}_M \sim p} [\|\text{VQ}(x, \mathcal{C}_M)\|^2]$ is radial is a consequence of Lemma 1. We can thus denote $M_{d,p,M}$ this function. Moreover, by Lemma 4 for any $x \in \mathbb{R}^d$, $\lambda \in (0, 1)$ and **almost surely** any codebook \mathcal{C}_M , we have that $\|\text{VQ}(\lambda x, \mathcal{C}_M)\|^2 \leq \|\text{VQ}(x, \mathcal{C}_M)\|^2$. (We apply Lemma 4 with $c_1 = \text{VQ}(x, \mathcal{C}_M)$ and $c_2 = \text{VQ}(\lambda x, \mathcal{C}_M)$). Consequently, for any $\lambda \in (0, 1)$, we have $\mathbb{E}_{\mathcal{C}_M \sim p} [\|\text{VQ}(\lambda x, \mathcal{C}_M)\|^2] \leq \mathbb{E}_{\mathcal{C}_M \sim p} [\|\text{VQ}(x, \mathcal{C}_M)\|^2]$, and $M_{d,p,M}$ is non-decreasing.

To prove the second statement, we use the following decomposition: for any $x \in \mathbb{R}^d$ such that $\|x\| = R$, and a.s. any \mathcal{C}_M :

$$\begin{aligned} \|\text{VQ}(x, \mathcal{C}_M)\|^2 &\leq 2 \|\text{VQ}(x, \mathcal{C}_M) - x\|^2 + 2 \|x\|^2 \\ &\leq 2 \|\text{VQ}(x, \mathcal{C}_M) - x\|^2 + 2R^2. \end{aligned}$$

681 Taking the expectation, and using again Lemma 1, we get

$$M_{d,p,M}(R) \leq 2\mathbb{E}_{\mathcal{C}_M \sim p} [\| \text{VQ}(Re_1, \mathcal{C}_M) - Re_1 \|^2] + 2R^2.$$

682 We finally prove that, for a given $x \in \mathbb{R}^d$, $\mathbb{E}_{\mathcal{C}_M \sim p} [\| \text{VQ}(x, \mathcal{C}_M) - x \|^2]$ is non-increasing with
 683 M . Indeed, writing $\mathcal{C}_M = [C_1, \dots, C_M] = \mathcal{C}_{M-1} \cup \{C_M\}$ (which amounts to define a **coupling**
 684 between \mathcal{C}_M and \mathcal{C}_{M-1}), we obtain

$$\begin{aligned} \| \text{VQ}(x, \mathcal{C}_M) - x \|^2 &= \min \left(\| \text{VQ}(x, \mathcal{C}_{M-1}) - x \|^2, \| C_M - x \|^2 \right) \\ &\leq \| \text{VQ}(x, \mathcal{C}_{M-1}) - x \|^2. \end{aligned}$$

685 □

686 In the next Lemma and the rest of the proof, to avoid the cumbersome notation $(r_M^p)^{-1}$ we omit the
 687 dependency on p in r_M^p and simply write r_M^{-1} . We now show that the quadratic error is uniformly
 688 decreasing on $B_2(0; R)$ at speed $M^{-2/d}$.

689 **Lemma 8.** *Assume A 3-A 4-A 5. Then, for any $R > 0$ and $M_0 \geq 1$ there exists a constant*
 690 $C_{M_0, R} < \infty$ *such that, for all $M \geq M_0$,*

$$\sup_{x \in B_2(0; R)} \mathbb{E}_{\mathcal{C}_M \sim p} \left[\| r_M^{-1}(\|x\|) \text{VQ}(x, \mathcal{C}_M) - x \|^2 \right] \leq C_{M_0, R} M^{-2/d}.$$

691 *Proof.* Using Lemma 1, the function $x \rightarrow \mathbb{E}_{\mathcal{C}_M \sim p} \left[\| r_M^{-1}(\|x\|) \text{VQ}(x, \mathcal{C}_M) - x \|^2 \right]$ is radial. For
 692 $x \in \mathbb{R}^d$, we get that

$$\begin{aligned} &\mathbb{E}_{\mathcal{C}_M \sim p} \left[\| r_M^{-1}(\|x\|) \text{VQ}(x, \mathcal{C}_M) - x \|^2 \right] \\ &\leq 2\mathbb{E}_{\mathcal{C}_M \sim p} \left[\| (r_M^{-1}(\|x\|) - 1) \text{VQ}(x, \mathcal{C}_M) \|^2 \right] + 2\mathbb{E}_{\mathcal{C}_M \sim p} \left[\| \text{VQ}(x, \mathcal{C}_M) - x \|^2 \right] \\ &\leq 2(r_M^{-1}(\|x\|) - 1)^2 \mathbb{E}_{\mathcal{C}_M \sim p} \left[\| \text{VQ}(x, \mathcal{C}_M) \|^2 \right] + 2\mathbb{E}_{\mathcal{C}_M \sim p} \left[\| \text{VQ}(x, \mathcal{C}_M) - x \|^2 \right]. \end{aligned} \quad (11)$$

693 Set $M_0 \geq 1$. First, by Lemma 7, for any $M \geq M_0$, we get

$$\sup_{x \in B_2(0; R)} \mathbb{E}_{\mathcal{C}_M \sim p} \left[\| \text{VQ}(x, \mathcal{C}_M) \|^2 \right] \leq C_{M_0, R, d, p} \quad (12)$$

694 Second, by Lemma 5, for any $x \in B_2(0; R)$, we have $0 \leq (r_M^{-1}(\|x\|) - 1) \leq (r_M^{-1}(R) - 1)$, and by
 695 the remark following Theorem 2, we get:

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

696 thus

$$(r_M^{-1}(R) - 1)^2 = |r_M(R) - 1|^2 / r_M^2(R) = O(M^{-2/d}). \quad (13)$$

697 Thirdly, Lemma 6 gives that

$$\sup_{x \in B_2(0; R)} \mathbb{E}_{\mathcal{C}_M \sim p} \left[\| \text{VQ}(x, \mathcal{C}_M) - x \|^2 \right] = \mathbb{E}_{\mathcal{C}_M \sim p} \left[\| \text{VQ}(Re_1, \mathcal{C}_M) - Re_1 \|^2 \right],$$

698 Using Theorem 2, we obtain

$$\lim_{M \rightarrow \infty} M^{2/d} \mathbb{E}_{\mathcal{C}_M \sim p} [\| \text{VQ}(Re_1, \mathcal{C}_M) - Re_1 \|^2] = C_d p_{\text{rad}}^{-2/d}(R). \quad (14)$$

699 Plugging (12)-(13) and (14) into (11), we obtain that

$$\sup_{x \in B_2(0; R)} \mathbb{E}_{\mathcal{C}_M \sim p} \left[\| r_M^{-1}(\|x\|) \text{VQ}(x, \mathcal{C}_M) - x \|^2 \right] = O(M^{-2/d}).$$

700 □

A random scalar quantizer is defined by a scalar output codebook, $\mathcal{O}_Q = [o_1, \dots, o_Q]$ assumed to be ordered $o_1 < \dots < o_Q$ and an external randomization, which may be taken without loss of generality as $U \sim \text{Unif}([0, 1])$; see Appendix B.1 for a construction. Here Q is the number of codewords and the number of bits required to encode the output vectors (the scalar quantizer rate) is $\log_2(Q)$. A scalar quantizer is said to be *uniform* if for all $i \in [Q - 1]$, $o_{i+1} - o_i = \delta$ for some $\delta > 0$. It is shown in Appendix B.1 that, for all $L > 0$, we may construct a uniform random scalar quantizer with Q codewords, satisfying for all $x \in [0, L]$,

$$\mathbb{E}_{U \sim \text{Unif}([0, 1])} [\text{SQ}(x, L, Q, U)] = x \quad (15)$$

$$\mathbb{E}_{U \sim \text{Unif}([0, 1])} [\{\text{SQ}(x, L, Q, U) - x\}^2] \leq L^2/4(Q - 1)^2. \quad (16)$$

We use this random scalar quantifier in the following proposition.

Proposition 1. Assume A 3-A 4-A 5. Let $R > 0$ and $M_0 \geq 1$. For $Q \geq 2$, denote by $\omega_{M, Q}(R) = \max_{x \in B_2(0; R)} \Omega_{M, Q}(x)$ where

$$\Omega_{M, Q}(x) = \mathbb{E}_{\mathcal{C}_M \sim p, U \sim \text{Unif}([0, 1])} \left[\left\| \text{SQ}(r_M^{-1}(\|x\|), r_M^{-1}(R), Q, U) \text{VQ}(x, \mathcal{C}_M) - x \right\|^2 \right].$$

Then, there exists a constant $C_{M_0}(R)$, such that for all $M \geq M_0$,

$$\omega_{M, Q}(R) \leq C_{M_0}(R) \{M^{-2/d} + Q^{-2}\}.$$

Proof. It follows from Lemma 1 that the function $x \rightarrow \Omega_{M, Q}(x)$ is radial. If A is a scalar random variable, \mathbf{B} is a random vector, A and \mathbf{B} are independent, $\mathbb{E}[A^2] < \infty$, and $\mathbb{E}[\|\mathbf{B}\|^2] < \infty$, then

$$\mathbb{E}[\|A\mathbf{B} - \mathbb{E}[A]\mathbb{E}[\mathbf{B}]\|^2] = \mathbb{E}[(A - \mathbb{E}[A])^2] \mathbb{E}[\|\mathbf{B}\|^2] + \{\mathbb{E}[A]\}^2 \mathbb{E}[\|\mathbf{B} - \mathbb{E}[\mathbf{B}]\|^2].$$

Setting $A \leftarrow \text{SQ}(r_M^{-1}(\|x\|))$, $\mathbf{B} \leftarrow \text{VQ}(x, \mathcal{C}_M)$, we get

$$\begin{aligned} \Omega_{M, Q}(x) &= \mathbb{E}_{\mathcal{C}_M \sim p} \left[\left\| r_M^{-1}(\|x\|) \text{VQ}(x, \mathcal{C}_M) - x \right\|^2 \right] \\ &+ \mathbb{E}_{\mathcal{C}_M \sim p} \left[\left\| \text{VQ}(x, \mathcal{C}_M) \right\|^2 \right] \mathbb{E}_{U \sim \text{Unif}([0, 1])} \left[\left\{ \text{SQ}(r_M^{-1}(\|x\|), r_M^{-1}(R), Q, U) - r_M^{-1}(\|x\|) \right\}^2 \right]. \end{aligned}$$

We now make the following observations:

1. The function $r_M^{-1}(\|x\|)$ is bounded on $B_2(0; R)$ by $r_M^{-1}(R)$ (see Lemma 5). Hence,

$$\mathbb{E}_{U \sim \text{Unif}([0, 1])} \left[\left\{ \text{SQ}(r_M^{-1}(\|x\|), r_M^{-1}(R), Q, U) - r_M^{-1}(\|x\|) \right\}^2 \right] \leq r_M^{-2}(R)/4(Q - 1)^2.$$

2. The second order moment $\mathbb{E}_{\mathcal{C}_M \sim p} \left[\left\| \text{VQ}(x, \mathcal{C}_M) \right\|^2 \right]$ is upper bounded by a constant independent of M on $B_2(0; R)$ by Lemma 7.

3. Using Lemma 8, we can upper bound the first term for any $M \geq M_0$

$$\mathbb{E}_{\mathcal{C}_M \sim p} \left[\left\| r_M^{-1}(\|x\|) \text{VQ}(x, \mathcal{C}_M) - x \right\|^2 \right] \leq C_{M_0, R} M^{-2/d}.$$

This concludes the proof. \square

We now give the proof of the main result Theorem 4.

Proof of Theorem 4. We consider the process described in DoStoVoQ. For $t \geq 0$, $k \in [K]$, we define $(\hat{b}_{k, t+1}^\ell)_{\ell \in [L]}$ such that $\hat{g}_{k, t+1} = \|g_{k, t+1}\| D^{-1/2} (\hat{b}_{k, t+1}^1, \dots, \hat{b}_{k, t+1}^L)$, where for any $\ell \in [L]$ we have $\hat{b}_{k, t+1}^\ell = \text{StoVoQ}(b_{k, t+1}^\ell, p, M, P, s_{k, t+1})$.

1. Conditional independence property. We observe that for all $k \in [K]$, $g_{k, t+1}$ is \mathcal{G}_{t+1} -measurable. Moreover, the seeds $(s_{k, t+1})_{k \in [K], t \geq 0}$ are independent, and there exists a functional ϕ such that for all $k \in [K]$, $\hat{g}_{k, t+1} = \phi(g_{k, t+1}, s_{k, t+1})$. We conclude that the compressed stochastic gradients $(\hat{g}_{k, t+1})_{k \in [K]}$ are mutually independent conditionally to \mathcal{G}_{t+1} .

729 **2. Unbiasedness.** In the sequel, we fix $t \geq 0$, $k \in [K]$. Regarding the bias, we use the fact that for
 730 any $\ell \in [L]$, $b_{k,t+1}^\ell$ is \mathcal{G}_{t+1} -measurable. Moreover, as $\hat{b}_{k,t+1}^\ell = \text{StoVoQ}(b_{k,t+1}^\ell, p, M, P, s_{k,t+1})$, we
 731 have that $\mathbb{E}_{\mathcal{C}_M \sim p}[\hat{b}_{k,t+1}^\ell | \mathcal{G}_{t+1}] = b_{k,t+1}^\ell$, using the fact that StoVoQ is unbiased.

732 Consequently, $\mathbb{E}_{\mathcal{C}_M \sim p} \left[\|g_{k,t+1}\| / \sqrt{D} \times (\hat{b}_{k,t+1}^1, \dots, \hat{b}_{k,t+1}^L)_{\ell \in [L]} \mid \mathcal{G}_{t+1} \right] = g_{k,t+1}$.

733 **3. Relative error bound.** We write:

$$\mathbb{E} [\|\hat{g}_{k,t+1} - g_{k,t+1}\|^2 \mid \mathcal{G}_{t+1}] = \frac{\|g_{k,t+1}\|^2}{D} \sum_{\ell \in [L]} \mathbb{E} [\|\hat{b}_{k,t+1}^\ell - b_{k,t+1}^\ell\|^2 \mid \mathcal{G}_{t+1}].$$

734 Remark that $\sum_{\ell \in [L]} \|b_{k,t+1}^\ell\|^2 = D$. Consequently, for all $\ell \in [L]$, $b_{k,t+1}^\ell \in B_2(0; \sqrt{D})$. Using
 735 Proposition 1, with $R = D$, we get:

$$\mathbb{E} [\|\hat{g}_{k,t+1} - g_{k,t+1}\|^2 \mid \mathcal{G}_{t+1}] = \frac{\omega_{M,Q}(R)}{D} \|g_{k,t+1}\|^2$$

736 which concludes the proof. □

737

738 B Scalar and vector Quantization

739 B.1 Unbiased random scalar quantization

740 A random scalar quantizer is a random map from the real line to a (scalar) codebook $\mathcal{O}_Q =$
 741 $\{o_1, \dots, o_Q\} \subset \mathbb{R}$ where $Q \geq 2$. It is assumed that $-\infty < o_1 < \dots < o_Q < \infty$. The reso-
 742 lution (or code rate) is $P = \log_2(Q)$ is the number of bits needed to uniquely specify a codeword.
 743 A scalar quantizer is said to be *uniform* if for all $i \in [Q-1]$, $o_{i+1} - o_i = \delta$, for some $\delta > 0$. Note
 744 that in such case $\delta = \{o_Q - o_1\} / (Q-1)$.

745 For $x \in \mathbb{R}$ and $u \in [0, 1]$, consider a function $\text{SQ}(x, \mathcal{O}_Q, u) \in \mathcal{O}_Q$. If $U \sim \text{Unif}([0, 1])$, then
 746 $\text{SQ}(x, \mathcal{O}_Q, U)$ is a random scalar quantizer. A random scalar quantizer is said to be *unbiased* if for
 747 all $x \in [o_1, \dots, o_Q]$, $\mathbb{E}_{U \sim \text{Unif}([0, 1])}[\text{SQ}(x, \mathcal{O}_Q, U)] = x$.

748 A simple way to construct an unbiased scalar quantizer goes as follows. We first compute the index
 749 $j(x) \in [Q]$ such that $x \in [o_{j(x)}, o_{j(x)+1})$. Note that $x = \lambda_{j(x)}^*(x) o_{j(x)} + (1 - \lambda_{j(x)}^*(x)) o_{j(x)+1}$
 750 where

$$\lambda_{j(x)}^*(x) = (x - o_{j(x)}) / (o_{j(x)+1} - o_{j(x)}) \in (0, 1].$$

751 For $u \in (0, 1]$, we set

$$\text{SQ}(x, \mathcal{O}_Q, u) = \mathbb{1}_{\{u \leq \lambda_{j(x)}^*(x)\}} o_{j(x)} + \mathbb{1}_{\{u > \lambda_{j(x)}^*(x)\}} o_{j(x)+1}.$$

752 Since $\mathbb{E}_{U \sim \text{Unif}([0, 1])}(U \leq \lambda_{j(x)}^*(x)) = \lambda_{j(x)}^*(x)$ the unbiasedness follows. It is easily seen that the
 753 distortion of a scalar quantizer is directly related to the diameter of the quantizer.

754 **Proposition 2.** For all $x \in [o_1, o_Q]$, it holds that

$$\mathbb{E}_{U \sim \text{Unif}([0, 1])}[\{\text{SQ}(x, \mathcal{O}_Q, U) - x\}^2] \leq (1/4) \sup_{i \in [Q-1]} \{o_{i+1} - o_i\}^2.$$

755 If the scalar quantizer is uniform,

$$\mathbb{E}_{U \sim \text{Unif}([0, 1])}[\{\text{SQ}(x, \mathcal{O}_Q, U) - x\}^2] \leq (1/4)(Q-1)^{-2} \{o_Q - o_1\}^2.$$

756 *Proof.* For all $x \in [o_1, o_Q]$, we get

$$|\text{SQ}(x, \mathcal{O}_Q, U) - x| \leq (1/2) \{o_{j(x)+1} - o_{j(x)}\}$$

757 The proof follows. □

758 Unbiased random scalar quantization is a special case of dual vector quantization, introduced in the
 759 next section.

760 B.2 Dual Vector Quantization

761 We introduce a new notion of vector quantization, called *dual quantization* (or *Delaunay quantiza-*
 762 *tion*). The principle of dual quantization is to map an \mathbb{R}^d -valued vector x onto a codebook \mathcal{C}_M using
 763 a random splitting operator $\text{Dual-VQ}(x, \mathcal{C}_M, U)$ such that, for all $x \in \text{ConvHull}(\mathcal{C}_M)$,

$$\mathbb{E}_{U \sim \text{Unif}([0,1])}[\text{Dual-VQ}(x, \mathcal{C}_M, U)] = x. \quad (17)$$

764 We stress that in this case the unbiasedness is not due to the use of a random codebook but makes use
 765 of an external randomization. In practice, a dual quantizer procedure amounts to define a probability
 766 distribution of \mathcal{C}_M , with weights $(\lambda_1^*(x), \dots, \lambda_M^*(x))$, $\lambda_i^*(x) \geq 0$, $\sum_{j=1}^M \lambda_j^*(x) = 1$. Set $\Lambda_0^*(x) = 0$
 767 and for $i \in [M]$, $\Lambda_i^*(x) = \sum_{j=1}^i \lambda_j^*(x)$. Note that $\Lambda_M^*(x) = 1$. If $u \in (\Lambda_{j-1}^*(x), \Lambda_j^*(x)]$, $j \in [M]$,
 768 we set $\text{Dual-VQ}(x, \mathcal{C}_M, u) = c_j$. In such that, for all $x \in \text{ConvHull}(\mathcal{C}_M)$, we get

$$\mathbb{E}_{U \sim \text{Unif}([0,1])}[\text{Dual-VQ}(x, \mathcal{C}_M, U)] = \sum_{i=1}^M \lambda_i^*(x) c_i = x.$$

769 The distortion of a dual quantizer is therefore given, for $x \in \mathbb{R}^d$, by

$$\mathbb{E}_{U \sim \text{Unif}([0,1])}[\|\text{Dual-VQ}(x, \mathcal{C}_M, U) - x\|^2] = \sum_{i=1}^M \lambda_i^*(x) \|x - c_i\|^2. \quad (18)$$

770 For $x \in \text{ConvHull}(\mathcal{C}_M)$, the probability distribution $(\lambda_1^*(x), \dots, \lambda_M^*(x))$ is obtained by solving the
 771 following convex optimization program:

$$(\lambda_1^*(x), \dots, \lambda_M^*(x)) = \underset{(\lambda_1, \dots, \lambda_M) \in \mathcal{S}(x, \mathcal{C}_M)}{\text{argmin}} \sum_{i=1}^M \lambda_i \|x - c_i\|^2, \quad (19)$$

772 where

$$\mathcal{S}(x, \mathcal{C}_M) = \left\{ (\lambda_1, \dots, \lambda_M) \in \mathbb{R}_+^M, \sum_{i=1}^M \lambda_i = 1, \sum_{i=1}^M \lambda_i c_i = x \right\}. \quad (20)$$

773 The support of $(\lambda_1^*(x), \dots, \lambda_M^*(x))$ is $M + 1$ at most. For a distribution q on \mathbb{R}^d , we define

$$\text{Dual-Dist}(q, \mathcal{C}_M) = \int q(x) \left\{ \sum_{i=1}^M \lambda_i^*(x) \|x - c_i\|^2 \right\} dx. \quad (21)$$

774 For a given input distribution q , an optimal codebook \mathcal{C}_M^* of cardinality M satisfies
 775 $\text{Dual-Dist}(q, \mathcal{C}_M^*) \leq \text{Dual-Dist}(q, \mathcal{C}_M)$ for all \mathcal{C}_M satisfying $|\mathcal{C}_M| = M$.

776 **Theorem 5** (Rates, see [27]). Asymptotic rate. Assume that the pdf q is compactly supported on
 777 \mathbb{R}^d -valued.

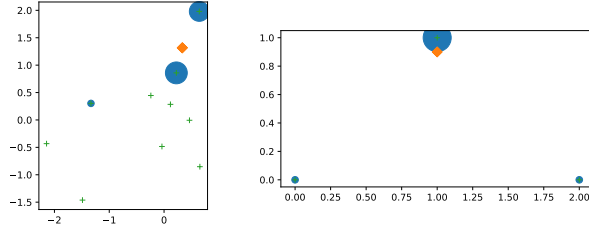
$$\lim_{M \rightarrow \infty} M^{\frac{2}{d}} \inf_{|\mathcal{C}_M| = M} \text{Dual-Dist}(q, \mathcal{C}_M) =: Q_2^D(q) = Q_2^D(\text{Unif}([0,1]^d)) \|q\|_{\frac{d}{d+2}}. \quad (22)$$

778 B.3 HSQ methods - see Dai et al. [8]

779 In this Section, we provide a detailed review of the two methods proposed by Dai et al. [8]. In
 780 Appendix B.3.1, we first discuss HSQ-Span and explain why it cannot compete with approaches
 781 based on Voronoi quantization. In Appendix B.3.2, we discuss HSQ-greed.

782 B.3.1 HSQ-Span

783 The first method, HSQ-Span, is unbiased but suffers form a large variance. Indeed, it relies on
 784 decomposing the vector $x \in \mathbb{R}^d$ as a **linear** combination of the codewords in \mathcal{C}_M , assuming that
 785 $\text{Span}\{c_i, i \in [M]\} = \mathbb{R}^d$ (a codebook satisfying this property is said to be *full-rank*). Because
 786 typically $M \gg d$, there are infinitely many solutions to the linear problem $\sum_{i=1}^M \alpha_i c_i = x$, i.e.
 787 $\mathcal{A}(x, \mathcal{C}_M) = \{(\alpha_1, \dots, \alpha_M) \in \mathbb{R}^M, \sum_{i=1}^M \alpha_i c_i = x\}$ is infinite. Note that contrary to the Dual
 788 quantization approach, we do not assume that $\alpha_i \geq 0$ for $i \in [M]$ or $\sum_{i \in [M]} \alpha_i = 1$. However, for



(a) $x \sim \mathcal{N}(0, I_2)$ and $M = 10$. (b) $x \sim \mathcal{N}(0, I_2)$ and $M = 3$.

Figure 3: Delaunay quantization for a vector x (orange diamond), for a given set of codewords (green +), and corresponding weights (area of the blue spheres). Remark that all but three points have a 0 probability of being picked, making the quadratic error much smaller than for HSQ-span.

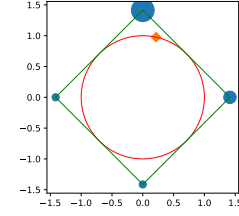


Figure 4: Cross-Polytope [10] is a particular case of Delaunay quantization. The codewords are the vertices of $B_1(0; \sqrt{d})$. A vector x (orange diamond) lying on the unit Ball $B_2(0; 1)$ (red circle) is decomposed with weights (area of the blue spheres) of codewords on the Ball of radius \sqrt{d} (green).

any $i \in [M]$, we pick the codeword c_i with probability $|\alpha_i|/\|\alpha\|_1$, and encode x as $\text{sign}(\alpha_i)\|\alpha\|_1 c_i$. In HSQ-Span, the **minimal norm** solution in $\mathcal{A}(x, \mathcal{C}_M)$ is chosen, i.e. solve

$$\alpha^*(x) := (\alpha_1^*(x), \dots, \alpha_M^*(x)) = \underset{(\alpha_1, \dots, \alpha_M) \in \mathcal{A}(x, \mathcal{C}_M)}{\text{argmin}} \sum_{i=1}^M \alpha_i^2, \quad (23)$$

The main advantage are that

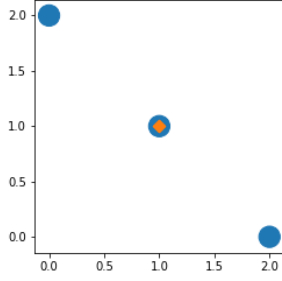
1. **Fast computation.** First, as $\alpha^*(x) = C^\dagger x$, where C^\dagger is the Moore-Penrose inverse of the codewords matrix $C = [c_1, \dots, c_M]$, provided a fixed codebook \mathcal{C}_M , it is possible to compute only once C^\dagger and to then obtain $\alpha^*(x)$ for any x by a simple matrix-vector product.
2. **Unbiased.** Second, this approach is unbiased. Its quadratic error thus linearly decays with the number of workers.

However (1) its variance is high and (2) does not decrease with M . Indeed, the minimal norm solution $\alpha^*(x)$ tends to put weight on **all** codewords. For example, we represent in Figure 6 the weights on each vector for 3 situations in dimension $d = 2$. Intuitively, the probability of selecting c_i is not a decreasing function $\|x - c_i\|^2$ (see e.g., Figure 6b), which results in the large variance; even if there exists i such that $c_i = x$, there is a non vanishing probability of selecting $c_j \neq c_i$ (Figure 6a). We illustrate the second point in Figure 5 which gives the evolution of the distortion for $d = 16$ w.r.t. M for $K \in \{1, 8\}$ workers. The error does not decrease.

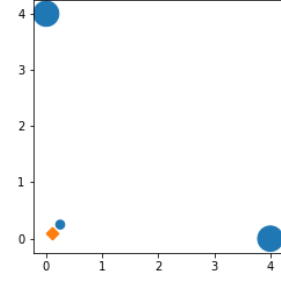
B.3.2 HSQ-greed

HSQ-greed is closed to DoStoVoQ: Dai et al. [8] still consider a full-rank codebook \mathcal{C}_M , and simply encode x by $\text{VQ}(x, \mathcal{C}_M)$. We list here the main differences to our approach:

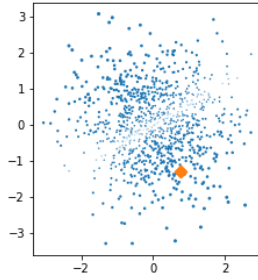
1. the same codebook is used during all iterations and on all workers. This makes it impossible or cumbersome to apply the convergence result developed in the federated learning literature, which require that the **compression on each workers are independent** (at least between iterations).
2. No assumption is made on the codebook distribution (apart from the fact that it is full-rank). The importance of unitary invariance is not mentioned. In practice, authors use an codebook generated by applying a k-means algorithm on a larger set of scaled Gaussian isotropic vectors. This pre-processing slightly improves the distribution of the codewords but is in practice of limited impact (see paragraph (e) in Subsection 2.2 in the main text.).



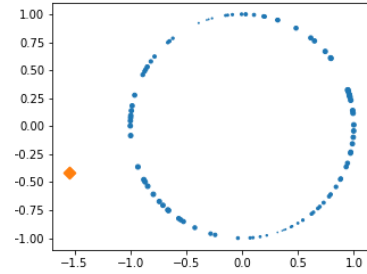
(a) $x_1 = (1, 1)$,
 $\mathcal{C}_M = [(2, 0); (0, 2); (1, 1)]$.
 $\alpha^*(x_1) \simeq (33\%, 33\%, 33\%)$.



(b) $x_2 = (0.1, 0.1)$,
 $\mathcal{C}_M = [(4, 0); (0, 4); (1/4, 1/4)]$
 $\alpha^*(x_2) \simeq (47\%, 47\%, 6\%)$.



(c) $x \sim \mathcal{N}(0, I_2)$ and $M = 1000$.



(d) $x \sim \mathcal{N}(0, I_2)$ and $M = 100$, $p = \mathcal{U}(\mathcal{S}_1(\mathbb{R}^2))$.

Figure 6: HSQ-Span: weights (size of the blue point) on each of the codewords of \mathcal{C}_M when decomposing x (orange diamond).

- 817 3. Codewords are chosen of norm 1. This means we also need to encode $\|x\|$ together
- 818 with $\text{VQ}(x, \mathcal{C}_M)$, which is typically done on using 6 bits per bucket.
- 819 4. The method is biased, so does not benefit from a large number of workers. No analysis of
- 820 the quadratic error is provided.

821 **Theoretical results.** Dai et al. [8] present a
 822 convergence result for HSQ-greed, namely in
 823 Lemma 3 and the subsequent Theorem 3. Note
 824 however that the proof of this result is **not** pro-
 825 vided in the paper³. Second, the guarantee pro-
 826 vided is almost vacuous. Indeed, authors rely
 827 on an alternative assumption⁴ on the **alignment**
 828 of the compressed value $\text{VQ}(x, \mathcal{C}_M)$ with x :

829 **Definition 5** (Compression with preserved
 830 alignment). *There exists $\alpha > 0$ such that*
 831 *for all $x \in \mathbb{R}^d$, we have $\langle \text{Comp}(x), x \rangle^2 \geq$*
 832 *$(1 - \alpha)\|x\|^2$.*

833 This assumption becomes stronger as $1 - \alpha$
 834 increases. However, Lemma 3 indicates that
 835 $1 - \alpha \geq \sigma_{\min}(C)/M$, with σ_{\min} the minimal

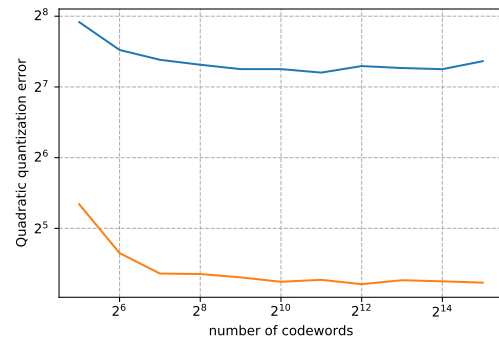


Figure 5: HSQ-Span: Distortion as a function of M (log-scale): $K = 1$ (blue) $K = 8$ (orange).

³The appendices of the paper were not available, neither on <https://arxiv.org/pdf/1911.04655.pdf>, nor on <https://paperswithcode.com/paper/hyper-sphere-quantization-communication>, on the 1st of June, 2021.

⁴See the work of Beznosikov et al. [5] for a discussion between the possible assumptions.

eigenvalue of the codebook matrix C . The bound guarantee thus worsens with M . A similar multiplicative factor $1/(1 - \alpha) \propto M$ appears in their convergence rate Theorem 3. We note that without any assumption on the codebook distribution, it seems difficult to obtain any result, as the worst case codebook that satisfies the full-rank assumption could be arbitrarily bad (typically a unique codeword perturbed by a tiny amount of noise $\mathcal{C}_M = [c_1 + \eta \epsilon_i]_{i \in [M]}$, with η very small).

B.4 Alignment under A 3, for StVoQ without debiasing function (new)

We can leverage our tight analysis of unitarily invariant distribution to obtain a result on the expected alignment between a vector $x \in \mathbb{R}^d$ and the output of StVoQ (without bias removal) applied on x . More precisely, we have the following lemma:

Lemma 9. *Assume Lemma 1, then :*

$$\mathbb{E}_{\mathcal{C}_M \sim p}[\langle x, \text{VQ}(x, \mathcal{C}_M) \rangle] \geq r_M^p(\|x\|)\|x\|^2.$$

We proved that on a ball of radius R , $r_M^p(\|x\|) \geq r_M^p(R)$ and that $r_M^p(R) = 1 - O(M^{-2d})$: in other words our guarantee **does** improve with M . This result is thus much stronger than the one of [8]. Note however that without debiasing, it is not possible to directly leverage the literature on federated learning: our result is only on the *expected* alignment and it would require novel proofs (relying on 1. expected alignment and 2. Bounded variance) to give a convergence result.

C Unitarily invariant random codebooks

We gather in this section the theoretical results on random codebooks distributed according to a unitarily invariant distribution. In Appendix C.1 we provide a proof of Theorem 2. In Appendix C.2, we show that random codebook are asymptotically optimal when the distribution of inputs is also unitarily invariant, provided that the codebook distribution is appropriately chosen. In Appendix C.3, we state an elementary lower bound. In Appendix C.4, we extend our results for spherical codebooks. **This extension is not mentioned in the main document.**

C.1 Proof of Theorem 2

Let $\{C_i\}_{i=1}^\infty$ be an i.i.d. sequence of random vector with pdf p w.r.t. the Lebesgue measure. For each $n \in \mathbb{N}$, denotes by $\mathcal{C}_M = \{C_1, \dots, C_M\}$ the associated sequence of codebook. Define for $k > 0$ and $\lambda > 0$, the pdf of the Weibull distribution (Weibull(k, λ)) with shape parameter k and scale parameter λ is given, for $x \geq 0$, by:

$$w_{k,\lambda}(x) = (k/\lambda)(x/\lambda)^{k-1}e^{-(x/\lambda)^k}. \quad (24)$$

The survival function of Weibull(k, λ) is denoted by $\bar{W}_{k,\lambda}(t) = e^{-(t/\lambda)^k}$. Denote by V_d the volume of the unit ball in \mathbb{R}^d ,

$$V_d = \pi^{d/2}/\Gamma(d/2 + 1). \quad (25)$$

We preface the proof by a lemma.

Lemma 10. *Assume A 3-A 4. Then, for every $x \in \mathbb{R}^d$ and $t \geq 0$,*

$$\lim_{M \rightarrow \infty} \mathbb{P}_{\mathcal{C}_M \sim p}(M^{1/d} \|\text{VQ}(x, \mathcal{C}_M) - x\| > t) = \bar{W}_{d, (V_d p(x))^{-1/d}}(t).$$

Proof. We get for $t \geq 0$,

$$\begin{aligned} \mathbb{P}_{\mathcal{C}_M \sim p}(M^{1/d} \min_{i=1:M} \|C_i - x\| \geq t) &= \left(1 - \mathbb{P}_{C_1 \sim p}(\|C_1 - x\| \leq tM^{-1/d})\right)^M \\ &= \left(1 - P(B(x; tM^{-1/d}))\right)^M \end{aligned}$$

where for $A \subset \mathbb{R}^d$ a Borel set, $P(A) = \int p(x) \mathbb{1}_A(x) dx$. It follows from the Lebesgue differentiation theorem that

$$P(B(x; tM^{-1/d})) \sim_{M \rightarrow \infty} p(x) \text{Leb}_d(B(x; tM^{-1/d})) = p(x) V_d t^d M^{-1}$$

870 where Leb_d denotes the Lebesgue measure. Hence, for any $t \geq 0$

$$\mathbb{P}_{\mathcal{C}_M \sim p} \left(M^{1/d} \min_{i=1:M} \|C_i - x\| \geq t \right) \xrightarrow{M \rightarrow \infty} \bar{W}_{d, \{V_d p(x)\}^{-1/d}}(t) = e^{-p(x)V_d t^d}.$$

871

□

872 *Proof of Theorem 2.* The proof relies on Lemma 10 and on the uniform integrability of the sequence

873 $(M^{1/d} \min_{i=1:M} \|C_i - x\|)^2$, $M \geq 1$. Let $R > 0$.

$$\begin{aligned} \mathbb{E}_{\mathcal{C}_M \sim p} \left[M^{2/d} \min_{i \in [M]} \|C_i - x\|^2 \mathbb{1}_{\{M^{2/d} \min_{i=1:M} \|C_i - x\|^2 \geq R\}} \right] \\ = \int_R^\infty \left\{ 1 - P(B(x; M^{-1/d} t^{1/2})) \right\}^M dt \end{aligned}$$

874 By Anderson's inequality (see [11]) $P(B(x; M^{-1/d} t^{1/2})) \geq P(B(0; M^{-1/d} t^{1/2}))$ so that

$$\int_R^{+\infty} \left\{ 1 - P(B(x; M^{-1/d} t^{1/2})) \right\} dt \leq \int_R^{+\infty} \underbrace{\left(1 - P(B(0; M^{-1/d} t^{1/2})) \right)^M}_{=: \phi_M(t)} dt.$$

875 Let $\rho \in (0, 1)$. Decompose $\phi_M(t) = A_M(t) + B_M(t)$ with

$$\begin{aligned} A_M(t) &= \left\{ 1 - P(B(0; M^{-1/d} t^{1/2})) \right\}^M \mathbb{1}_{\{M^{-1/d} t^{1/2} > t^{\rho/2}\}} \\ B_M(t) &= \exp \left(-M \left[\frac{P(B(0; M^{-1/d} t^{1/2}))}{\text{Leb}_d(B(0; M^{-1/d} t^{1/2}))} V_d M^{-1/d} t^{d/2} \right] \right) \mathbb{1}_{\{M^{-1/d} t^{1/2} \leq t^{\rho/2}\}}. \end{aligned}$$

876 Note that

$$A_M(t) \leq 1 - P(B(0; t^{\rho/2})).$$

877 Now let $\delta > 0$. We upper-bound $B_M(t)$ as follows

$$\begin{aligned} B_M(t) &\leq \exp \left(- \inf_{s \in (0, \delta]} \frac{P(B(0; s))}{\text{Leb}_d(B(0; s))} V_d t^{d/2} \right) + \exp \left(- \inf_{s \in (\delta, t^{\rho/2}]} \frac{P(B(0; s))}{\text{Leb}_d(B(0; s))} V_d t^{d/2} \right) \\ &\leq \exp \left(- \inf_{s \in (0, \delta]} \frac{P(B(0; s))}{\text{Leb}_d(B(0; s))} V_d t^{d/2} \right) + \exp \left(-P(B(0; \delta)) t^{d(1-\rho)/2} \right). \end{aligned}$$

878 Let us denote $B_1(t)$ and $B_2(t)$ the two terms on the right hand side of the previous equation. Note

879 that $B_i(t)$, $i = 1, 2$ do not depend on M . Now let us show that these functions are integrable.

$$\int_0^{+\infty} A_M(t) dt \leq \int_0^{+\infty} \left\{ 1 - P(B(0; t^{\rho/2})) \right\} dt = \int_0^{+\infty} \mathbb{P}(\|C_1\| > t^{\rho/2}) dt = \mathbb{E}[\|C_1\|^{2/\rho}] < +\infty.$$

880 For the next two terms we use the elementary inequality $\inf_{s \in (0, \delta]} \frac{P(B(0; s))}{\lambda_d(B(0; s))} \geq m_\delta$ so that

$$\int_0^{+\infty} B_1(t) dt \leq \int_0^{+\infty} e^{-V_d m_\delta t^{d/2}} dt < +\infty$$

881 and

$$\int_0^{+\infty} B_2(t) dt \leq \int_0^{+\infty} e^{-m_\delta V_d t^{d(1-\rho)/2}} dt < +\infty.$$

882 Consequently, $\lim_{R \rightarrow +\infty} \sup_M \int_R^{+\infty} \phi_M(t) dt = 0$ which ensures uniform integrability.

883 We conclude by using that the second moment of a Weibull(k, λ) is given by $\lambda^2 \Gamma(1 + 2/k)$ and that

884 $V_d = \pi^{d/2} / \Gamma(1 + d/2)$. □

885 C.2 Proof of Theorem 3

886 The proof of Theorem 3 follows almost immediately from Theorem 2 using the classical Zador
887 theorem, stated for completeness below (see [14] for a proof of Zador theorem, which has a long
888 history).

889 **Theorem 6** (Zador's Theorem).

890 • Assume that $\int \|x\|^{r+\delta} p(x) dx < \infty$, for some $\delta > 0$. Then,

$$\lim_{M \rightarrow \infty} M^{2/d} \text{Dist}(p, \mathcal{C}_M^*) = Q_2(p) = Q_2([0, 1]^d) \left(\int p^{d/(d+2)}(x) dx \right)^{(d+2)/d}, \quad (26)$$

891 where $Q_2(p)$ is the quantization coefficient of the distribution p and $Q_2([0, 1]^d)$ that of the uni-
892 form distribution over the unit hypercube, $\text{Unif}([0, 1]^d)$. If the distribution p is standard normal
893 $\mathcal{N}(0, \mathbf{I}_d)$, then $Q_2(\mathcal{N}(0, \mathbf{I}_d)) \sim_{d \rightarrow \infty} d$.

894 • There exists a universal constant $C_{d,r+\delta} \in (0, \infty)$ such that, for any pdf p on \mathbb{R}^d

$$\text{Dist}(p, \mathcal{C}_M^*) \leq C_{d,r+\delta} \sigma_{r+\delta}^r M^{-r/d}$$

895 where $\sigma_{r+\delta}(p) =: \inf_{a \in \mathbb{R}^d} \left(\int |x|^{r+\delta} |x| p(x) dx \right)^{\frac{1}{r+\delta}}$.

896 C.3 An elementary lower-bound

897 The Hölder inequality with negative exponents (see [16, p. 191]) shows that for $0 <$
898 $r < 1$ and $s \in \mathbb{R}$ such that $r^{-1} + s^{-1} = 1$ (hence $s < 0$), $\int p^{-2/d}(x) q(x) dx \geq$
899 $\left\{ \int p^{-2s/d}(x) dx \right\}^{1/s} \left\{ \int q^r(x) dx \right\}^{1/r}$. Setting $s = -d/2$ and $r = 2/(d+2)$, we get that
900 $C(q, p, d) \geq \|q\|_{d/(d+2)}$.

901 C.4 Asymptotic distortion of a random quantizer on the unit sphere S_{d-1}

902 We now consider random codebooks on the unit hypersphere $S_{d-1} = \{x \in \mathbb{R}^d, \|x\| = 1\}$. We
903 compute the distortion of a codebook distributed uniformly on the unit-sphere as a function of the
904 number of codewords M and of the ambient dimension d . Denote by σ_{d-1} the uniform distribution
905 on S_{d-1} . Denote

$$\kappa_d = \left(\frac{2\sqrt{\pi} \Gamma((d+1)/2)}{\Gamma(d/2)} \right)^{1/(d-1)} \quad (27)$$

906 **Theorem 7.** Assume $d \geq 2$ and assume that the codewords $\{C_n\}_{n \geq 1}$ are i.i.d. uniformly distributed
907 on the unit hyper-sphere S_{d-1} of \mathbb{R}^d . For every $x \in S_{d-1}$, and $t \geq 0$

$$\lim_{n \rightarrow \infty} \mathbb{P}_{\mathcal{C}_M \sim \sigma_{d-1}} (M^{1/(d-1)} \|VQ(x, \mathcal{C}_M) - x\| \geq t) = \bar{W}_{d-1, \kappa_d}(t), \quad .$$

908 Furthermore,

$$\lim_{M \rightarrow \infty} M^{2/(d-1)} \mathbb{E}_{\mathcal{C}_M \sim \sigma_{d-1}} [\|VQ(x, \mathcal{C}_M) - x\|^2] = \kappa_d^2 \Gamma(1 + 2/(d-1)).$$

909 *Proof.* Since the uniform distribution over S_{d-1} is unitary invariant, we get for all $x \in S_{d-1}$,

$$\mathbb{P}_{C_1 \sim \sigma_{d-1}} (C_1 \in B(x; r)) = \frac{\sigma_{d-1}(B(x; r) \cap S_{d-1})}{\sigma_{d-1}(S_{d-1})} \sim_{r \rightarrow 0^+} \frac{\text{Leb}_{d-1}(B(0; r))}{\sigma_{d-1}(S_{d-1})}$$

910 and, using that $\sigma_{d-1}(S_{d-1}) = 2\pi^{d/2}/\Gamma(d/2)$ we get

$$\frac{\text{Leb}_{d-1}(B(0; r))}{\sigma_{d-1}(S_{d-1})} = \frac{V_{d-1} r^{d-1}}{\sigma_{d-1}(S_{d-1})} = \frac{\pi^{\frac{d-1}{2}} r^{d-1}}{\Gamma(\frac{d-1}{2} + 1) 2\pi^{\frac{d}{2}}} = \frac{\Gamma(\frac{d}{2})}{\Gamma(\frac{d+1}{2})} \frac{r^{d-1}}{2\sqrt{\pi}} = (r/\kappa_d)^{d-1}.$$

911 Taking $r = M^{-1/(d-1)} t$ yields that

$$\begin{aligned} \mathbb{P}_{\mathcal{C}_M \sim \sigma_{d-1}} (M^{-1/(d-1)} \|VQ(x, \mathcal{C}_M) - x\| > t) &= \left(1 - \mathbb{P}_{C_1 \sim \sigma_{d-1}} (C_1 \in B(0; M^{-1/(d-1)} t)) \right)^M \\ &\longrightarrow_{M \rightarrow \infty} e^{-(t/\kappa_d)^{d-1}}. \end{aligned}$$

912 This completes the first part of the proof. For the second part of the proof, it is required to check
913 the uniform integrability which follows from the fact that the above equivalence (\sim) also holds as a
914 lower bound inequality. \square

915 D Algorithmic extensions

916 D.1 Spherical codebooks

917 In this section, we describe a spherical version of StoVoQ and DoStoVoQ. Beyond the obvious
 918 change from the codeword distribution from Gaussian to uniform on the sphere, a key modification
 919 stems from the fact that each quantized vector has norm 1: the debiasing function does not depend
 920 on $\|x\|$, but only on the number of codewords M . Consequently, the bias correction does not need
 921 to be transmitted and can be directly performed on the central server.

922 On the other hand, the norm of each bucket has to be transmitted: the vector quantization is applied
 923 to the *shape*, i.e. the unitary vector $x/\|x\|$. We use a scalar quantizer for the norm, typically over
 924 4-6 bits. For completeness, the codes of those two algorithms are given in Algorithms 3 and 4.

Algorithm 3: Spherical-StoVoQ

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

Output: Codeword index \mathbf{i}_c , value \mathbf{i}_r

```

925 1 Sample  $\mathcal{C}_M \sim \sigma_{d-1}$  with seed  $s$ ; /* sample codebook with uniform distribution  $\sigma_{d-1}$  on the sphere */
2  $c_l = \text{VQ}(x/\|x\|, \mathcal{C}_M)$ ; /* quantize (select a codeword in spherical codebook  $\mathcal{C}_M$ ) */
3  $\mathbf{i}_{c_l} \leftarrow \text{index of } c_l$ ; /* get index of codeword */
4  $\mathbf{i}_r = \text{SQ}(\|x\|)$ ; /* quantize  $r$  on  $P$  bits */

```

Algorithm 4: Spherical-DoStoVoQ over T iterations

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

Output: $(\theta_t)_{t \geq 0}$

```

1 for  $t = 1, \dots, T$  do
2    $w_0$  sends  $\theta_{t-1}$  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  $\theta_{t-1}$ ;
5     Split  $g_{k,t}$  on  $[b_{k,t}^1, \dots, b_{k,t}^L]$ ;
926 6   for  $\ell = 1, \dots, L$  (in parallel) do
7      $(\mathbf{i}_c^{t,k,\ell}, \mathbf{i}_r^{t,k,\ell}) = \text{Spherical-StoVoQ}(b_{k,t}^\ell, p, M, P, s_{k,t})$ 
8   end
9   Send  $(\mathbf{i}_c^{t,k,\ell}, \mathbf{i}_r^{t,k,\ell})_{\ell \in [L]}$  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

```

927 D.2 Extension to DoStoVoQ-DIANA and DoStoVoQ-VR-DIANA

928 In this subsection, we provide the adaptations of the DoStoVoQ algorithm to algorithms designed to
 929 handle heterogeneous workers, and for which the best complexities are achieved, namely DIANA [25]
 930 and VR-DIANA [17]. Those algorithms are based on the fundamental idea: relying on control variates
 931 $(h_{k,t})_{k \in [K], t \geq 0}$, updated at each iteration, that converge (in the convex case), for each worker k , to
 932 $\nabla f_k(\theta^*)$. Instead of compressing $g_{k,t}$, the algorithm compresses the difference between the actual
 933 gradient and the control variate $g_{k,t} - h_{k,t}$. The impact of those control variates (often referred to as
 934 *memory*) is to mitigate the discrepancy between workers' gradients that stems from the heterogeneity
 935 of the data-distribution between different workers. As explained in Appendix E it is particularly
 936 relevant to reduce this discrepancy to maximize the impact of the multiple workers. The same idea
 937 can be incorporated within a variance reduced algorithm, we here focus on SVRG [18] (extension to
 938 SAGA [9] or other variants is straightforward). To incorporate variance reduction to the algorithm,
 939 we further assume that each f_k is a finite sum $\frac{1}{S} \sum_{s \in [S]} f_{k,s}$. Algorithms DoStoVoQ-DIANA and
 940 DoStoVoQ-DIANA-SVRG are provided in respectively Algorithms 5 and 6.

Algorithm 5: DoStoVoQ-DIANA over T iterations . Lines specific to the Diana approach are highlighted in **blue**

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

```

1 Set  $h_{k,0} = 0$  for all  $k \in [K]$  (or alternatively  $h_{0,k} = \nabla f_k(\theta_0)$ );
2 for  $t = 1, \dots, T$  do
3    $w_0$  sends  $\theta_{t-1}$  and different seeds  $s_{k,t}$  to each  $w_k$ ;
4   for  $k = 1, \dots, K$  do
5     Compute local gradient  $g_{k,t}$  at  $\theta_{t-1}$ ;
6     Set  $\Delta_{k,t} = g_{k,t} - h_{k,t}$ ;
7     Split  $\Delta_{k,t} \times \sqrt{D}/\|\Delta_{k,t}\|$  on  $[\delta_{k,t}^1, \dots, \delta_{k,t}^L]$ ;
8     for  $\ell = 1, \dots, L$  (in parallel) do
9        $(\mathbf{i}_c^{t,k,\ell}, \mathbf{i}_r^{t,k,\ell}) = \text{StoVoQ}(\delta_{k,t}^\ell, p, M, P, s_{k,t})$ 
10    end
11    Reconstruct  $(\hat{\Delta}_{k,t})_{k \in K}$ ;
12    Update memory:  $h_{k,t+1} = h_{k,t} + \alpha \hat{\Delta}_{k,t}$ ;
13    Send  $(\|\Delta_{k,t}\|, (\mathbf{i}_c^{t,k,\ell}, \mathbf{i}_r^{t,k,\ell})_{\ell \in [L]})$  to  $w_0$ ;
14  end
15  On the central node  $w_0$ ;
16  Reconstruct  $(\hat{\Delta}_{k,t})_{k \in K}$ ;
17  Update:  $\theta_t = \theta_{t-1} - \gamma_t(\bar{h}_t + \frac{1}{K} \sum_{k=1}^K \hat{\Delta}_{k,t})$ ;
18  Update averaged memory :  $\bar{h}_{t+1} (:= \frac{1}{K} \sum_{k \in [K]} h_{k,t}) = \bar{h}_t + \frac{\alpha}{K} \sum_{k \in [K]} \hat{\Delta}_{k,t}$ ;
19 end

```

Algorithm 6: DoStoVoQ-DIANA-SVRG over T iterations . Lines specific to the variance reduction approach are highlighted in **green**

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

```

1 Set  $h_{k,0} = 0$  for all  $k \in [K]$  (or alternatively  $h_{0,k} = \nabla f_k(\theta_0)$ );
2 for  $t = 1, \dots, T$  do
3   Sample  $u_t \sim \mathcal{B}(S^{-1})$ ;
4    $w_0$  sends  $\theta_{t-1}, u_t$  and different seeds  $s_{k,t}$  to each  $w_k$ ;
5   for  $k = 1, \dots, K$  do
6     if  $u_t = 1$  then
7       Set  $\eta_{k,s,t} = \theta_t$  for all  $s \in [S]$ ;
8       Sample  $s_{k,t} \sim \text{Unif}[S]$ ;
9       Set  $\mu_{t,k} = S^{-1} \sum_{s \in S} \nabla f_{k,s}(\eta_{k,s,t})$ ;
10      Set  $g_{k,t} = \nabla f_{k,s_{k,t}}(\theta_{t-1}) - \nabla f_{k,s_{k,t}}(\eta_{k,s_{k,t},t}) + \mu_{t,k}$ ;
11      Set  $\Delta_{k,t} = g_{k,t} - h_{k,t}$ ;
12      Split  $\Delta_{k,t} \times \sqrt{D}/\|\Delta_{k,t}\|$  on  $[\delta_{k,t}^1, \dots, \delta_{k,t}^L]$ ;
13      for  $\ell = 1, \dots, L$  (in parallel) do
14         $(\mathbf{i}_c^{t,k,\ell}, \mathbf{i}_r^{t,k,\ell}) = \text{StoVoQ}(\delta_{k,t}^\ell, p, M, P, s_{k,t})$ 
15      end
16      Reconstruct  $(\hat{\Delta}_{k,t})_{k \in K}$ ;
17      Update memory:  $h_{k,t+1} = h_{k,t} + \alpha \hat{\Delta}_{k,t}$ ;
18      Send  $(\|\Delta_{k,t}\|, (\mathbf{i}_c^{t,k,\ell}, \mathbf{i}_r^{t,k,\ell})_{\ell \in [L]})$  to  $w_0$ ;
19    end
20    On the central node  $w_0$ ;
21    Reconstruct  $(\hat{\Delta}_{k,t})_{k \in K}$ ;
22    Update:  $\theta_t = \theta_{t-1} - \gamma_t(\bar{h}_t + \frac{1}{K} \sum_{k=1}^K \hat{\Delta}_{k,t})$ ;
23    Update averaged memory :  $\bar{h}_{t+1} (:= \frac{1}{K} \sum_{k \in [K]} h_{k,t}) = \bar{h}_t + \frac{\alpha}{K} \sum_{k \in [K]} \hat{\Delta}_{k,t}$ ;
24  end

```

E Additional experiments

In this Section, we compare by Monte Carlo the distortions achieved by different compression schemes for 3 types of input x . For a given (random) compressor generically denoted $Q(\cdot)$ and $x \in \mathbb{R}^d$, we decompose $Q(x) = Q_{\parallel}(x) + Q_{\perp}(x)$, where $Q_{\parallel}(x) = \|x\|^{-2}xx^{\top}Q(x)$. With these notations, $Q_{\parallel}(x)$ and $Q_{\perp}(x)$ are the components of the quantization error which are colinear and orthogonal to x , respectively. By construction, $\|x - Q(x)\|^2 = \|x - Q_{\parallel}(x)\|^2 + \|Q_{\perp}(x)\|^2$. The distortion is computed for $K = 1$ and $K \in \{8, 20\}$ workers (depending on the experiments). We compare 10 compression schemes, corresponding to 7 algorithms (some with several variants): the signed algorithm (Sign) (see Definition 1), Top-H with $H = 2$ (see Definition 2), Rand-H with $H = 2$ (see Definition 3), HSQ-Span (see Appendix B.3.1) with $M = 2^{10}$ and a 6 bits scalar quantizer for the norm, HSQ-greed (see appendix B.3.2) with $M = 2^{10}$ and a 6 bits scalar quantizer for the norm, Polytope (see Appendix B.2) with and without quantization of the norm, three variants of StoVoQ with a Gaussian random codebook with $M = 2^{13}$ and $p = \mathcal{N}(0, (1 + 2/d)\mathbf{I}_d)$: GRVQ which is StoVoQ without the radial debiasing step, StoVoQ without quantization of r_M^p , and StoVoQ with an unbiased scalar quantization of $(r_M^p)^{-1}$ over $P = 3$ bits (strictly speaking, only this last column corresponds to the algorithm StoVoQ, the two previous versions have been added to assess the influence of the debiasing by $\{r_M^p\}^{-1}$ and the quantization of $\{r_M^p\}^{-1}$).

We compare those algorithms over three tasks:

1. **Task 1:** Compression of a random vector from a standard Gaussian input distribution in dimension $d = 16$. We compare $K = 1$ to $K = 20$. Results are given in Appendix E.1.
2. **Task 2:** Compression of “real” gradients, extracted from a training performed with a VGG16 on CIFAR10 with SGD, extracted at epoch 10, on which a pre-processing similar to DoStoVoQ is applied. The minibatch gradients on a given layer are divided into buckets of dimension $d = 16$. A normalisation is applied to sets of $L = 32$ buckets (the normalisation for the blocks of $d \times L = 512$ coefficients are scalar quantized with a high-resolution scalar quantizer and sent to the parameter server). Results are given in Appendix E.2.1. We compare the performance with 1 and 8 workers, when *all workers compress the same gradient*. The goal of this task is to assess the impact of the actual distribution of the normalised minibatch gradients values w.r.t. a Gaussian distribution.
3. **Task 3:** Compression of “real” gradients, with multiple workers, each worker compresses a different minibatch stochastic gradient, computed at the same parameter (as described in Subsection 4.2): this is the most practical setting, and we explain the resulting trade-offs, especially in terms of the distribution of stochastic gradient noise (see [29]) and the inhomogeneity between workers. We perform this task on (i) the same setting as for task 2, and (ii) the gradients from the LS experiment introduced in Subsection 4.1. Results are given in Appendix E.2.2.

E.1 Distortion for Gaussian input

Setup: We here compare all methods on a Gaussian input $x \sim \mathcal{N}(0, \mathbf{I}_d)$ for $d = 16$. Monte Carlo is performed over 10^4 repetitions. Standard deviation is negligible.

Observations. Results are provided in Table 5. We make the following observations:

1. We first observe that in the single worker case, Sign, Top-2, HSQ-greed and StoVoQ-GRVQ achieve a global error or respectively 6.4, 8.7, 9.1 and 6.8. (These errors are obtained by summing the radial and orthogonal numbers). StoVoQ achieves an error of 11 which is slightly higher, Rand-2, Polytope, HSQ-span suffer a much higher errors of 110, 121, 147. This confirms the theoretical predictions.
2. We observe in practice here the fundamental differences between biased / unbiased compression methods and also methods that ensure the independence of the compression on each individual worker: while all biased methods do not benefit from the multiplicity of workers, for unbiased and independent compression, the distortion for $K = 1$ is divided by K . Here, the quadratic errors, both radial and orthogonal, are reduced by a factor 20. Overall, over 20 workers, the error obtained by StoVoQ, with debiasing and scalar quantization is 0.5. This is by far the best method in terms of global distortion for 20 workers.

Table 5: **Task 1:** Distortion for Gaussian inputs

Method	Sign	Top-2	Rand-2	Polytope	
Variant	norm-quant.				
$K = 1$	1.0 5.4	4.8 3.9	12 98	5.8 115	5.8 115
$K = 20$	1.0 5.4	4.7 3.8	0.6 4.8	0.3 5.6	0.3 5.6

Method	HSQ-span	HSQ-greed	StoVoQ		
Variant	norm-quant.	norm-quant.	GRVQ	Unbiased	Unbiased+quant.
$K = 1$	3.8 143	1.3 7.8	1.8 5.0	0.5 10.5	0.5 10.5
$K = 20$	0.2 7.0	1.3 7.5	1.7 0.25	0.03 0.5	0.03 0.5

- 995 3. **StoVoQ-GRVQ vs StoVoQ-Unbiased.** For StoVoQ, the application of debiasing increases the non-
996 radial quantization distortion, by a factor of nearly 2 (from 5 to 10), while simultaneously re-
997 ducing the radial distortion, from 2 to 0.5. This increase is unavoidable to obtain the unbiased
998 character, that is necessary to reduce the error beyond 1. Indeed, it is important to remark the
999 radial bias for StoVoQ-GRVQ and HSQ is not negligible (1.3 and 1.8 respectively): in fact , this
1000 radial distortion is also of order $M^{-2/d}$ thus using an even larger codebook would not reduce it
1001 significantly.
- 1002 4. **StoVoQ-Unbiased vs StoVoQ-Unbiased+Scalar-Quantization.** We observe that the impact of the
1003 scalar quantization is negligible here, which indicates that the impact of the scalar quantization
1004 of the norm is limited.
- 1005 5. **HSQ vs StoVoQ-GRVQ:** These two methods are somehow similar for a single worker: HSQ relies
1006 on a gain-shape decomposition with the a scalar quantization of the norm and a vector quanti-
1007 zation of the normalized vector using spherical codebooks whereas GRVQ uses random Gaussian
1008 codebooks with a variance matched to the input variance. We observe that overall StoVoQ-GRVQ
1009 slightly outperforms HSQ for $K = 1$. This is in favor of using Gaussian codebooks.

1010 E.2 Distortion for neural networks gradients

1011 E.2.1 Task 2: Impact of the distribution

1012 **Setup.** We compare the distortion for $K = 1$
1013 and 8, on stochastic gradients sampled from
1014 the training of a VGG16, with SGD, at epoch
1015 10. The gradients are partitioned into blocks of
1016 size 2^9 ; then each block is scaled and split into
1017 buckets of dimension $d = 16$. Those buckets
1018 are then compressed using each of the possible
1019 methods in dimension 16. The results presented
1020 are obtained using 1000 stochastic gradient .

1021 The main objective is to compare the impact
1022 of the distribution of the gradients on the dis-
1023 tortion of the different compressors. For ex-
1024 ample, if the stochastic gradient noise is heavy
1025 tailed (leading equivalently to sparse gradients),
1026 methods relying on sparsification, e.g., *Top-2*,
1027 is expected to perform significantly better than
1028 Gaussian random codebook (recall that the opti-
1029 mality result Theorem 3 assumes that the distribu-
1030 tion of the codewords matches the distribution
1031 of the inputs; the choice of a Gaussian distribu-
1032 tion for the codewords implicitly assumes that the
1033 distribution of the gradients is approximately Gaussian).

1031 For $K = 8$, we assume that each workers compresses the same gradient : we compute the error of
1032 $K^{-1} \sum_{k \in [K]} \hat{g}_{t,k} - g_t$, where $\hat{g}_{t,k}$ stands for the output of the k -th compressor on g_t .

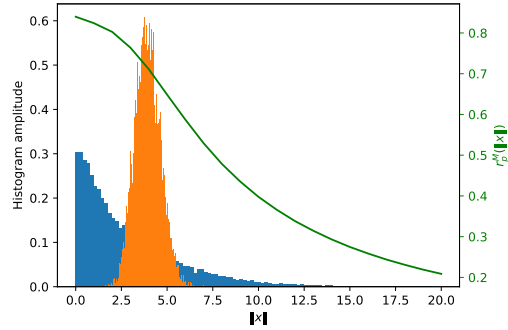


Figure 7: Histograms of the VGG16 gradient buckets (blue), of Gaussian vectors (orange), and the radial bias for the associated dimension $d = 16$ (green).

Table 6: **Task 2:** empirical distortion from a sample of gradients sampled from a VGG-16 at epoch 10, (same gradients on each worker).

Method	Sign	Top-2	Rand-2	Polytope	
Variant	norm-quant.				
$K = 1$	0.46 0.35	0.35 0.16	0.86 4.8	0.50 7.0	0.50 7.1
$K = 8$	0.46 0.35	0.35 0.16	0.15 0.61	0.07 0.9	0.07 0.9

Method	HSQ-span	HSQ-greed	StoVoQ		
Variant	norm-quant.	norm-quant.	GRVQ	Unbiased	Unbiased+quant.
$K = 1$	0.24 7.5	0.09 0.5	0.05 0.2	0.02 0.36	0.02 0.4
$K = 8$	0.07 0.9	0.09 0.4	0.04 0.03	0.002 0.05	0.003 0.05

As shown in Appendix A.5, the distortion of StoVoQ is a non-decreasing function of the norm of the vector to be compressed. In Figure 7, we represent simultaneously the histogram of the bucket norms, and the histogram of the norms of the Gaussian vectors used in **Task 1**. This suggests a departure from the Gaussian distribution for the stochastic gradient noise.

Observations Results are provided in Table 6. We highlight the following points.

1. Again, the unbiased version of StoVoQ achieves the best distortion.
2. Even though the distribution of the norms is very different from the norm of the Gaussian vectors (as illustrated in Figure 7), the distortion of StoVoQ is not severely impaired. Especially, the error for StoVoQ for a single worker is 0.4 vs 0.7 for Top-2, while for Gaussian inputs it was 11 vs 8.7 for Top-2.

E.2.2 Task 3: Signal-Noise ratio on the various gradients

Setup: We now consider that each worker computes and compresses a different stochastic gradient. More precisely, we collect samples of the stochastic gradients during an epoch: $[g_{t,1}^\top, \dots, g_{t,K}^\top]$, where $g_{t,k}$ is computed by the worker k on distinct minibatch of size b (all the gradients are $\{g_{t,k}\}_{k=1}^K$ are evaluated for the same value of the parameters). The compressed version is denoted $\{\hat{g}_{t,k}\}_{k=1}^K$.

In the homogeneous setting, for all $k \in [K]$, $g_{t,k}(\theta_t) = \nabla F(\theta_t) + \epsilon_{t,k}$, with $(\epsilon_{t,k})_{t,k}$ is the stochastic gradient noise $\mathbb{E}[\epsilon_{t,k} | \mathcal{F}_{t-1}] = 0$, where \mathcal{F}_{t-1} collects the past observations.

The central node averages the quantized stochastic gradient sent by the workers: $\tilde{g}_t := K^{-1} \sum_{k=1}^K \hat{g}_{t,k}$. We report in Tables 4 and 8 the normalized averaged error defined as

$$T^{-1} \sum_{t \in [T]} \frac{\|\frac{1}{K} \sum_{k=1}^K \hat{g}_{t,k} - g_{t,k}\|^2}{\|\frac{1}{K} \sum_{k=1}^K g_{t,k}\|^2}. \quad (28)$$

We here discuss in which settings we expect the multiple workers to improve w.r.t. a single worker. More precisely, we show that the impact of enforcing unbiased independent compression for the different workers increases with the "dependence" of stochastic gradients. Consider the following two cases. **Example 1: (large noise, low correlation between stochastic gradients)** each worker computes a stochastic gradient that is nearly independent of the other workers. The error made is **not** reduced by the multiplicity of workers. **Example 2: (low or no noise, strong consensus between stochastic gradients)** if each worker computes the same gradient, we recover **task 2**. The variance reduction obtained by using multiple workers and independent compressors is proportional to the number of workers. More generally, this is true when the noise is small w.r.t. the gradient of the function.

This signal/noise ratio fundamentally impacts the performance of algorithms using compressions operators: in **example 2**, it is crucial to use unbiased and independent workers, while in **example 1**, it is more important to reduce the distortion for a single worker.

Table 7: **Task 3:** normalized distortion for a mini-batch of size 4096 of a VGG-16 at epoch 10.

Method	Sign	Top-2	Rand-2	Polytope	
Variant	norm-quant.				
$K = 1$	0.3 0.2	0.5 0.2	0.5 6.2	0.2 7.3	0.2 7.3
$K = 8$	0.3 0.1	0.5 0.1	0.09 1.8	0.06 2.0	0.06 2.0

Method	HSQ-span	HSQ-greed	StoVoQ		
Variant	norm-quant.	norm-quant.	GRVQ	Unbiased	Unbiased+quant.
$K = 1$	0.2 8.5	0.09 0.5	0.06 0.2	0.02 0.4	0.02 0.4
$K = 8$	0.09 2.3	0.09 0.4	0.1 0.07	0.01 0.1	0.01 0.1

Many factors impact this “consensus” between workers: first of all the mini-batch size. The noise variance is inversely proportional to b : as b increases, each stochastic gradient becomes closer to $\nabla F(\theta)$. More generally, all variance reduction techniques tend to increase the “consensus”. On the other hand, heterogeneity between workers increases the discrepancy between gradients (but memory techniques as in DoStoVoQ-DIANA mitigate this discrepancy). Finally, performing several steps [36], as in Local-SGD also has a similar impact of averaging the noise over several iterations, and thus increases the consensus.

We thus evaluate all algorithms on two tasks:

1. First, on gradients extracted from the LSR task: in this task, data is distributed on all workers, that compute a batch gradient. The gradients obviously depend on the workers (each worker has access to a different subset of the data), but because the workers are homogeneous, these gradients have a strong consensus. We give the results in Table 8. We observe that the reduction by a factor 8 is preserved when using $K = 8$. This explains why our method outperforms HSQ in Figure 2.
2. Second, on gradients from the VGG16 trained with SGD on CIFAR. On this task, the noise level is much higher and the consensus much weaker. This is expected in very high dimensional models and non convex objective (roughly speaking, the gradients on different workers nearly point in random descent directions). We thus do not see any strong effect of the number of workers on the distortion for $b \leq 512$. Increasing further the batch size, to $b = 4096$, we recover the gain of multiple workers. Results are given in Table 7. The distortion is twice smaller with StoVoQ-unbiased than with any other method. While a batch of 4096 is very high, very large batch were used in a successful training of CIFAR and IMAGENET in Lin et al. [24]. More generally, when communication cost is a major concern, increasing the batch size and the number of local iterations is natural, to increase the quality of updates transmitted.

Table 8: **Task 3:** normalized distortion for LSR (see Section 4).

Method	Sign	Top-2	Rand-2	Polytope	
Variant	norm-quant.				
$K = 1$	0.05 0.3	0.4 0.2	0.6 6.3	0.3 7.3	0.3 7.3
$K = 8$	0.04 0.09	0.4 0.08	0.1 1.3	0.07 1.4	0.07 1.4

Method	HSQ-span	HSQ-greed	StoVoQ		
Variant	norm-quant.	norm-quant.	GRVQ	Unbiased	Unbiased+quant.
$K = 1$	0.3 9.4	0.09 0.5	0.1 0.3	0.03 0.6	0.03 0.6
$K = 8$	0.08 1.9	0.09 0.1	0.1 0.06	0.008 0.1	0.008 0.1