

Supplementary Document: Quantized Compressive Sampling of Stochastic Gradients for Efficient Communication in Distributed Training

Afshin Abdi, Faramarz Fekri*

School of Electrical and Computer Engineering
Georgia Institute of Technology, Atlanta, GA
{abdi,fekri}@gatech.edu

1 Dithered Quantization

It is well-known that the error in ordinary quantization, especially when the number of quantization levels is low, depends on the input signal and is not necessarily uniformly distributed. In *Dithered Quantization*, a (pseudo-)random signal called dither is added to the input signal prior to quantization. Adding this controlled perturbation can cause the statistical behavior of the quantization error to be more desirable (Schuchman 1964; Gray and Stockham 1993; Gray and Neuhoff 1998).

Let $Q(\cdot)$ be an M -level uniform quantizer with quantization step size of Δ , i.e., $Q(v) = \Delta \lfloor v/\Delta \rfloor$ and the output range of $Q(\cdot)$ is $\{-M, \dots, 0, \dots, M\}$.¹ The dithered quantizer is defined as follows;

Definition (Dithered Quantization). For an input signal x , let u be a dither signal, independent of x . The dithered quantization of x is defined as $\tilde{x} = Q(x + u) - u$.

Remark 1. To transmit the dithered quantization of x , it is sufficient to send the index of the quantization bin that $x + u$ resides in, i.e., $\lfloor (x + u)/\Delta \rfloor$. The receiver reproduces the (pseudo-)random sequence u using the same random number generator algorithm and seed number as the sender. It is then subtracted from $Q(x + u)$ to reconstruct \tilde{x} .

Theorem I ((Schuchman 1964)). *If 1) the quantizer does not overload, i.e., $|x + u| \leq \frac{M\Delta}{2}$ for all input signals x and dither u , and 2) The characteristic function of the dither signal, defined as $M_u(j\nu) = \mathbb{E}_u[e^{j\nu u}]$, satisfies $M_u(j\frac{2\pi l}{\Delta}) = 0$ for all $l \neq 0$, then the quantization error $e = x - \tilde{x}$ is uniform over $(-\Delta/2, \Delta/2]$ and it is independent of the signal x .*

It is common to consider $\mathcal{U}(-\Delta/2, \Delta/2)$ as the distribution of the random dither signal which can be easily verified that it satisfies the conditions of Thm. I.

*This material is based upon work supported by the National Science Foundation under Grant Number CPS-1837369 and by Sony Inc. under the Sony Faculty Research Award. Copyright © 2020, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

¹Throughout the paper, we assume that all quantizers are centered around 0 (with the exception of sign-based quantization). This is the case also for ternary (Wen et al. 2017) and stochastic quantization (Alistarh et al. 2017).

In some cases, the receiver may not be able to reproduce the dither signal to subtract from $Q(x + u)$. Hence, quantization is simply defined as $\tilde{x}_h = Q(x + u)$. We refer to this approach as the *half-dithered quantization* as the dither signal is applied only to the quantization, not the reconstruction of x . In this case, the quantization error is not necessarily independent of the signal, however by an appropriate choice of the dither signal, the moments of the quantization error will be independent (Gray and Stockham 1993). For example, if the dither signal u is the sum of k independent random variables, each having uniform distribution $\mathcal{U}(-\Delta/2, \Delta/2)$, then the k -th moment of the quantization error, $\epsilon = x - \tilde{x}_h$, would be independent of the signal; $\mathbb{E}[\epsilon^k | x] = \mathbb{E}[\epsilon^k]$.

Remark 2 (1-Bit Dithered Quantization). Note that the output range of the dithered quantization is $\{-M, \dots, +M\}$. Hence, each value is represented by minimum of $\log_2(1 + 2M)$ bits (without applying any compression to the quantized sequence). Reducing the number of bits to only 1-bit while keeping the desired properties of the dithered quantizers can potentially reduce the transmission bits by almost 50% (from at least $\log_2 3 \approx 1.58$ bits to 1 bit).

Without loss of generality, assume that $|x| \leq 1/2$. We propose the following dithered 1-bit quantization:

$$q = \text{sign}(x + u) := \begin{cases} +1 & \text{if } u + x > 0 \\ -1 & \text{o.w.} \end{cases}, \quad (1)$$

where $u \sim \mathcal{U}(-1/2, 1/2)$ is the random dither signal. The dequantized value is then given by

$$\tilde{x} = q - u. \quad (2)$$

It is straightforward to show that this 1-bit dithered quantizer is unbiased and the quantization noise is uniformly distributed and independent of x ;

$$\mathbb{E}[\tilde{x} - x] = 0, \quad \text{Var}[\tilde{x} - x] = \frac{1}{12}. \quad (3)$$

Relationship with Ternary and Stochastic Quantizations
Without loss of generality, assume that the vector x is normalized such that $|x_i| \leq 1$. Although the reconstruction of quantized values in our method is different from those in TernGrad and QSGD, we show that these quantizers can be considered as a special case of the half-dithered quantizer.

M -level Stochastic Quantization in (Alistarh et al. 2017) is defined as

$$Q^{(s)}(x_i) = \begin{cases} \text{sign}(x_i) \frac{l}{M} & \text{with prob. } 1 - d(x_i) \\ \text{sign}(x_i) \frac{l+1}{M} & \text{with prob. } d(x_i) \end{cases} \quad (4)$$

where l is the quantization bin that $|x_i|$ resides in, i.e., $|x_i| \in [l/M, (l+1)/M]$ and $d(x_i) := M|x_i| - l$. The ternary quantizer of (Wen et al. 2017) can be considered as a special case of stochastic quantizer with $M = 1$.

Lemma. *Stochastic quantization is the same as $(2M + 1)$ -level half-dithered quantizer with step-size $\Delta = \frac{1}{M}$ and uniform dither $u \sim \mathcal{U}(-\frac{1}{2M}, \frac{1}{2M})$.*

Proof. Let $Q(\cdot)$ be a $2M + 1$ -level quantizer with step size $\Delta = 1/M$. Let $u \sim \mathcal{U}[-\Delta/2, \Delta/2]$ be the dither signal. Let $0 \leq x \leq 1$ be an arbitrary number. Assume that $l/M \leq x < (l+1)/M$ and define $d = x - l/M$. Note that $0 \leq d < \Delta$ and

$$\begin{aligned} P\left(Q(x+u) = \frac{l}{M}\right) &= P\left(|x+u - l/M| \leq \frac{\Delta}{2}\right) \\ &= P\left(u \leq \frac{\Delta}{2} - d\right) = 1 - \frac{d}{\Delta} = 1 - Md. \end{aligned}$$

Similarly, $P(Q(x+u) = (l+1)/M) = Md$. Comparing with stochastic quantizer, we see that they both assign the quantization points with the same probability. The case $x < 0$ can be verified similarly. ■

In other words, stochastic quantizer adds a uniformly distributed dither to the input signal before quantization. However, the receiver *does not* subtract the dither from the quantized value. Therefore, the quantization error is not independent of the signal. It can be easily verified that although the quantization is unbiased, $\mathbb{E}[x - Q^{(s)}(x)] = 0$, its variance depends on the value of the input signal and varies in $[0, 1/4M^2]$, depending on the input signal x ;

$$\mathbb{E}\left[(Q^{(s)}(x) - x)_i^2\right] = \frac{d(x_i)(1 - d(x_i))}{M^2}.$$

On the other hand, the variance of the dithered quantization noise would be uniformly $1/12M^2$, independent of x . For example, if x is uniformly distributed over $[-1, 1]$, the average quantization variance of the stochastic quantizer would be $1/6M^2$, twice the variance of the dithered quantization.

2 Properties of the Mixing Matrix

Lemma II. *The random mixing matrix \mathbf{T} has the following properties:*

$$\mathbf{T}\mathbf{T}^\top = \frac{n}{k}\mathbf{I}, \quad \mathbb{E}[\mathbf{T}^\top\mathbf{T}] = \mathbf{I}. \quad (5)$$

Proof. Recall that $\mathbf{T} = \frac{1}{\sqrt{k}}\mathbf{H}\mathbf{R}$. Hence,

$$\mathbf{T}\mathbf{T}^\top = \frac{1}{k}\mathbf{H}\mathbf{R}\mathbf{R}^\top\mathbf{H}^\top \stackrel{(a)}{=} \frac{1}{k}\mathbf{H}\mathbf{H}^\top \stackrel{(b)}{=} \frac{n}{k}\mathbf{I},$$

where (a) is due to the fact that $\mathbf{R} = \text{diag}(\mathbf{r})$ and $r_i^2 = 1$, and (b) is a result of \mathbf{H} being any k rows of Hadamard matrix \mathbf{H}_n satisfying $\mathbf{H}_n\mathbf{H}_n^\top = n\mathbf{I}$.

For the second property,

$$\mathbb{E}[\mathbf{T}^\top\mathbf{T}] = \frac{1}{k}\mathbb{E}[\mathbf{R}\mathbf{H}\mathbf{H}^\top\mathbf{R}^\top].$$

Now, consider an arbitrary (i, j) -th element,

$$[\mathbf{R}^\top\mathbf{H}^\top\mathbf{H}\mathbf{R}]_{i,j} = r_i r_j [\mathbf{H}^\top\mathbf{H}]_{i,j}.$$

On the other hand, $\mathbb{E}[r_i r_j] = 1$ if $i = j$ and 0 if $i \neq j$. Moreover, since $H_{i,l} = \pm 1$, $[\mathbf{H}^\top\mathbf{H}]_{i,i} = \sum_{l=1}^k (H_{i,l})^2 = k$. Therefore,

$$\mathbb{E}[[\mathbf{R}^\top\mathbf{H}^\top\mathbf{H}\mathbf{R}]_{i,j}] = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases}$$

■

Note on Remark 1 of Paper Let M_n be an $n \times n$ random Rademacher matrix, i.e., each entry is an i.i.d. random variable $P(+1) = P(-1) = 0.5$. (Bourgain, Vu, and Wood 2010) has shown that

$$P(M_n \text{ is singular}) \leq \left(\frac{1}{\sqrt{2}} + o(1)\right)^n. \quad (6)$$

Now, for a constant c , consider collection of $m = cn/k$ of the QCS mixing matrix, $\mathbf{T}_i = \mathbf{H}\mathbf{R}_i$ for $i = 1, \dots, m$, where $\mathbf{R}_i = \text{diag}(\mathbf{r}_i)$ and \mathbf{r}_i 's are i.i.d. Rademacher random variables. Note that each of \mathbf{T}_i 's are full row-rank and the probability of \mathbf{T}_i 's not covering the entire \mathbb{R}^n would be less than the probability of M_n being singular.

3 Unbiased QCS

Proof of Thm. 2

For a fixed \mathbf{g} , we note that the randomness in $\hat{\mathbf{g}}$ stems from the random mixing matrix \mathbf{T} and dither signal \mathbf{u} , which are independent of each other and \mathbf{g} . Moreover, the quantization noise of \mathbf{v} can be written as

$$\hat{\mathbf{v}} = \mathbf{v} - \zeta\boldsymbol{\varepsilon} = \mathbf{T}\mathbf{g} - \zeta\boldsymbol{\varepsilon}, \quad (7)$$

where as a result of Thm. I, $\boldsymbol{\varepsilon}$ is an independent random variable and $\boldsymbol{\varepsilon} \sim \mathcal{U}(-1/2, 1/2)$.² Therefore, the quantization noise can be decomposed as

$$\mathbf{e} = \mathbf{g} - \hat{\mathbf{g}} = \overbrace{(\mathbf{I} - \mathbf{T}^\top\mathbf{T})\mathbf{g}}^{\mathbf{e}_g} + \overbrace{\zeta\mathbf{T}^\top\boldsymbol{\varepsilon}}^{\mathbf{e}_d}. \quad (8)$$

Unbiasedness.

$$\mathbb{E}[\mathbf{e}_g] = (\mathbf{I} - \mathbb{E}[\mathbf{T}^\top\mathbf{T}])\mathbf{g} \stackrel{(a)}{=} \mathbf{0},$$

$$\mathbb{E}[\mathbf{e}_d] = \mathbb{E}[\zeta\mathbf{T}^\top\boldsymbol{\varepsilon}] \stackrel{(b)}{=} \mathbb{E}[\zeta\mathbf{T}^\top] \mathbb{E}[\boldsymbol{\varepsilon}] = \mathbf{0},$$

where (a) is due to $\mathbb{E}[\mathbf{T}^\top\mathbf{T}] = \mathbf{I}$ (Lemma II) and (b) is because of independence of $\boldsymbol{\varepsilon}$ from $\zeta\mathbf{T}$ and $\mathbb{E}[\boldsymbol{\varepsilon}] = \mathbf{0}$. This proves the unbiasedness of QCS.

²Note that this is not the case for ordinary quantization or stochastic quantization of (Alistarh et al. 2017).

Variance. Note that since \mathbf{e}_g is a function of only \mathbf{T} and \mathbf{g} , it is independent of ε and

$$\mathbb{E}[\mathbf{e}_g^\top \mathbf{e}_g] = \mathbb{E}[\varepsilon \mathbf{e}_g^\top \mathbf{T}^\top] \mathbb{E}[\varepsilon] = 0.$$

Therefore,

$$\mathbb{E}[\|\mathbf{e}\|_2^2] = \mathbb{E}[\|\mathbf{e}_g\|_2^2] + \mathbb{E}[\|\mathbf{e}_d\|_2^2]$$

For an arbitrary \mathbf{g} , note that

$$\begin{aligned} \mathbb{E}[\|\mathbf{T}\mathbf{g}\|_2^2] &= \mathbf{g}^\top \mathbb{E}[\mathbf{T}^\top \mathbf{T}] \mathbf{g} = \mathbf{g}^\top \mathbf{I} \mathbf{g} = \|\mathbf{g}\|_2^2 \\ \mathbb{E}[\|\mathbf{T}^\top \mathbf{T} \mathbf{g}\|_2^2] &= \mathbb{E}[\mathbf{g}^\top \mathbf{T}^\top \mathbf{T} \mathbf{T}^\top \mathbf{T} \mathbf{g}] \\ &= \frac{n}{k} \mathbb{E}[\mathbf{g}^\top \mathbf{T}^\top \mathbf{T} \mathbf{g}] = \frac{n}{k} \|\mathbf{g}\|_2^2. \end{aligned}$$

Therefore,

$$\begin{aligned} \mathbb{E}[\|\mathbf{e}_g\|_2^2] &= \mathbb{E}[\|(\mathbf{I} - \mathbf{T}^\top \mathbf{T})\mathbf{g}\|_2^2] \\ &= \mathbb{E}[\|\mathbf{T}^\top \mathbf{T} \mathbf{g}\|_2^2] + \|\mathbf{g}\|_2^2 - 2 \mathbb{E}[\|\mathbf{T}\mathbf{g}\|_2^2] \\ &= \left(\frac{n}{k} - 1\right) \|\mathbf{g}\|_2^2. \end{aligned}$$

On the other hand,

$$\begin{aligned} \mathbb{E}[\|\mathbf{e}_d\|_2^2] &= \mathbb{E}[\|\zeta \mathbf{T}^\top \varepsilon\|_2^2] = \mathbb{E}_{\mathbf{T}}[\mathbb{E}_{\varepsilon}[\|\zeta \mathbf{T}^\top \varepsilon\|_2^2 | \mathbf{T}]] \\ &= \mathbb{E}_{\mathbf{T}}[\zeta^2 \mathbb{E}_{\varepsilon}[\|\mathbf{T}^\top \varepsilon\|_2^2 | \mathbf{T}]] \stackrel{(c)}{=} \mathbb{E}_{\mathbf{T}}\left[\zeta^2 \frac{\|\mathbf{T}\|_F^2}{12}\right] \\ &\stackrel{(d)}{=} \frac{n}{12Q^2} \mathbb{E}_{\mathbf{T}}[\|\mathbf{T}\mathbf{g}\|_\infty^2], \end{aligned}$$

where (c) is due to ε being i.i.d. $\mathcal{U}(-1/2, 1/2)$ and (d) is because of $\|\mathbf{T}\|_F^2 = n$ and definition of $\zeta = \|\mathbf{T}\mathbf{g}\|_\infty / Q$.

To bound $\mathbb{E}_{\mathbf{T}}[\|\mathbf{T}\mathbf{g}\|_\infty^2]$ we need the following lemma.

Lemma III. Let $\mathbf{a} \in \mathbb{R}^n$ be fixed and \mathbf{r} be an i.i.d. Rademacher random vector. Then for all $0 \leq \lambda < 1/(2\|\mathbf{a}\|_2^2)$,

$$\mathbb{E}_{\mathbf{r}}[e^{\lambda(\mathbf{a}^\top \mathbf{r})^2}] \leq \frac{1}{\sqrt{1 - 2\lambda\|\mathbf{a}\|_2^2}} \quad (9)$$

Proof. Let $\omega \sim \mathcal{N}(0, 1)$ be an independent normal random variable. Note that for a fixed \mathbf{r} ,

$$e^{\lambda(\mathbf{a}^\top \mathbf{r})^2} = \mathbb{E}_{\omega} \left[\exp \left(\sqrt{2\lambda}(\mathbf{a}^\top \mathbf{r})\omega \right) \right].$$

Therefore,

$$\begin{aligned} \mathbb{E}_{\mathbf{r}}[e^{\lambda(\mathbf{a}^\top \mathbf{r})^2}] &= \mathbb{E}_{\mathbf{r}} \left[\mathbb{E}_{\omega} \left[\exp \left(\sqrt{2\lambda}(\mathbf{a}^\top \mathbf{r})\omega \right) \mid \mathbf{r} \right] \right] \\ &= \mathbb{E}_{\omega} \left[\mathbb{E}_{\mathbf{r}} \left[\exp \left(\sqrt{2\lambda}(\mathbf{a}^\top \mathbf{r})\omega \right) \mid \omega \right] \right] \\ &\stackrel{(e)}{=} \mathbb{E}_{\omega} \left[\prod_{i=1}^n \mathbb{E}_{r_i} \left[\exp \left(\sqrt{2\lambda}\omega a_i r_i \right) \mid \omega \right] \right] \\ &\stackrel{(f)}{\leq} \mathbb{E}_{\omega} \left[\prod_{i=1}^n \exp \left(\lambda \omega^2 a_i^2 \right) \right] \\ &= \mathbb{E}_{\omega} \left[\exp \left(\lambda \|\mathbf{a}\|_2^2 \omega^2 \right) \right] = \frac{1}{\sqrt{1 - 2\lambda\|\mathbf{a}\|_2^2}} \end{aligned}$$

where (e) is because of independence of r_i 's, (f) is from Hoeffding's lemma, \blacksquare

Using the above lemma, for $\mathbf{v} = \mathbf{T}\mathbf{g}$ and arbitrary $\lambda > 0$,

$$\begin{aligned} \mathbb{E}[\|\mathbf{v}\|_\infty^2] &\leq \frac{1}{\lambda} \mathbb{E} \left[\log \left(\sum_{i=1}^k \exp(\lambda v_i^2) \right) \right] \\ &\leq \frac{1}{\lambda} \log \mathbb{E} \left[\sum_{i=1}^k \exp(\lambda v_i^2) \right] \\ &= \frac{1}{\lambda} \log \left(\sum_{i=1}^k \mathbb{E}[\exp(\lambda v_i^2)] \right) \end{aligned}$$

Note that for an arbitrary fixed i , $\zeta_j := r_j H_{i,j}$ would be i.i.d. Rademacher random variables and $v_i = \sum_j (\frac{g_j}{\sqrt{k}}) \zeta_j$.

Therefore, by Lemma III

$$\mathbb{E}[\exp(\lambda v_i^2)] \leq \frac{1}{\sqrt{1 - 2\lambda\|\mathbf{g}\|_2^2/k}},$$

and

$$\mathbb{E}[\|\mathbf{v}\|_\infty^2] \leq \frac{1}{\lambda} \log \left(\frac{k}{\sqrt{1 - 2\lambda\|\mathbf{g}\|_2^2/k}} \right).$$

For $k \geq 2$, let $\lambda = \frac{k}{2\|\mathbf{g}\|_2^2} (1 - k^{-\delta})$. Therefore,

$$\begin{aligned} \frac{1}{\lambda} \log \left(\frac{k}{\sqrt{1 - 2\lambda\|\mathbf{g}\|_2^2/k}} \right) &= \frac{2\|\mathbf{g}\|_2^2}{k} \frac{1}{1 - k^{-\delta}} \log(k^{1+\delta/2}) \\ &= \|\mathbf{g}\|_2^2 \frac{\log(k)}{k} \frac{2 + \delta}{1 - k^{-\delta}}. \end{aligned} \quad (10)$$

Setting $\delta = 1$, results in the bound

$$\mathbb{E}[\|\mathbf{v}\|_\infty^2] \leq 3\|\mathbf{g}\|_2^2 \frac{\log(k)}{k-1}. \quad (11)$$

Summarizing the above results for $k \geq 2$, we have

$$\mathbb{E}[\|\hat{\mathbf{g}} - \mathbf{g}\|_2^2] \leq \left(\frac{n}{k} - 1 + \frac{n}{4Q^2} \frac{\log(k)}{k-1} \right) \|\mathbf{g}\|_2^2. \quad (12)$$

For $k = 1$, note that since v is a scalar, by the definition of the used dithered quantizer sending magnitude of v and its sign results in $\hat{v} = v$ and $\mathbf{e}_d = 0$. Therefore,

$$\mathbb{E}[\|\hat{\mathbf{g}} - \mathbf{g}\|_2^2] = (n-1)\|\mathbf{g}\|_2^2. \quad (13)$$

Lower Bound on the Compression Gain of Sparsification

Recall that the variance of the unbiased sparsification is given by ((5) of (Wangni et al. 2018))

$$\mathbb{E}[\|\mathbf{g} - \hat{\mathbf{g}}\|_2^2] = \sum_{i=1}^n \frac{g_i^2}{p_i} - \|\mathbf{g}\|_2^2, \quad (14)$$

where $0 \leq p_i \leq 1$ and $\sum_i p_i = s$ controls the average transmission rate, i.e., the average compression gain is $\frac{n}{s}$. Finding the optimum value for p_i 's for general \mathbf{g} is not straightforward and hard to analyze. Hence, we assume that $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and compute the expected value of the sparsification variance over \mathbf{g} .

Using the inequality $\frac{a^2}{x} + \frac{b^2}{y} \geq \frac{(a+b)^2}{x+y}$ for arbitrary positive numbers (Titu's Lemma), the variance of error can be lower bounded as

$$\mathbb{E}[\|\mathbf{g} - \hat{\mathbf{g}}\|_2^2] \geq \frac{(\sum_i |g_i|)^2}{\sum_i p_i} - \|\mathbf{g}\|_2^2 = \frac{\|\mathbf{g}\|_1^2}{s} - \|\mathbf{g}\|_2^2.$$

Therefore,

$$\gamma = \frac{\mathbb{E}[\|\mathbf{g} - \hat{\mathbf{g}}\|_2^2]}{\|\mathbf{g}\|_2^2} \geq \frac{1}{s} \frac{\|\mathbf{g}\|_1^2}{\|\mathbf{g}\|_2^2} - 1.$$

On the other hand, for Gaussian random vector, $\mathbb{E}[\|\mathbf{g}\|_1^2] = \sum_{i,j} \mathbb{E}[|g_i| |g_j|] = n(1 + (n-1)\frac{2}{\pi})$, $\mathbb{E}[\|\mathbf{g}\|_2^2] = n$, and $\mathbb{E}[\frac{\|\mathbf{g}\|_1^2}{\|\mathbf{g}\|_2^2}]$ is close to that ratio. Especially, for large n , $\|\mathbf{g}\|_1^2$ and $\|\mathbf{g}\|_2^2$ are concentrated around their mean. Therefore,

$$\bar{\gamma} = \mathbb{E}[\gamma] \gtrsim \frac{1}{s} \left(\frac{2}{\pi} n + (1 - \frac{2}{\pi}) \right) - 1 \geq \frac{2}{\pi} \frac{n}{s} - 1.$$

4 MMSE QCS

Proof of Lemma 3

For $\hat{\mathbf{g}} = \alpha \mathbf{T}^\top \hat{\mathbf{v}}$, using the same argument as for the unbiased QCS, it can be easily shown that

$$\begin{aligned} \mathbb{E}[\|\mathbf{g} - \hat{\mathbf{g}}\|_2^2] &= \mathbb{E}[\|\mathbf{g} - \alpha \mathbf{T}^\top \mathbf{T} \mathbf{g} + \alpha \varsigma \mathbf{T}^\top \boldsymbol{\varepsilon}\|_2^2] \\ &= \|\mathbf{g}\|_2^2 (1 - 2\alpha + \alpha^2 \frac{n}{k}) + \alpha^2 \frac{n}{12Q^2} \mathbb{E}[\|\mathbf{T} \mathbf{g}\|_\infty^2] \\ &\leq \|\mathbf{g}\|_2^2 (1 - 2\alpha + \alpha^2 \frac{n}{k} + \alpha^2 \frac{n}{4Q^2} \frac{\log(k)}{k-1}) \\ &= \|\mathbf{g}\|_2^2 (1 - 2\alpha + \alpha^2 (\gamma + 1)), \end{aligned}$$

where $\gamma = \frac{n}{k} + \frac{n}{4Q^2} \frac{\log(k)}{k-1}$ for $k \geq 2$, and $\gamma = n - 1$ for $k = 1$.

Minimizing the upper bound of the error results in

$$\alpha = \frac{1}{1 + \gamma},$$

and by substituting α , the minimum mean squared error is given by

$$\mathbb{E}[\|\mathbf{g} - \hat{\mathbf{g}}\|_2^2] \leq \|\mathbf{g}\|_2^2 \frac{\gamma}{1 + \gamma}$$

5 Weighted Error Feedback

In weighted error feedback, instead of directly compressing the SGs, we add a weighted residue of compression from previous step to the SG and then compress it.

$$\mathbf{z}_t \leftarrow \mathbf{g}_t + \beta \mathbf{r}_t \quad (15a)$$

$$\hat{\mathbf{z}}_t \leftarrow \text{COMPRESS}(\mathbf{z}_t) \quad (15b)$$

$$\mathbf{e}_t \leftarrow \mathbf{z}_t - \hat{\mathbf{z}}_t \quad (15c)$$

$$\mathbf{r}_{t+1} \leftarrow (1 - \beta) \mathbf{r}_t + \mathbf{e}_t \quad (15d)$$

Proof of Lemma 4

Recall that $\mathbf{e}_t = \mathbf{z}_t - \hat{\mathbf{z}}_t = (\mathbf{I} - \alpha \mathbf{T}_t^\top \mathbf{T}_t) \mathbf{z}_t - \alpha \varsigma \mathbf{T}_t^\top \boldsymbol{\varepsilon}$, where $\alpha = 1$ for Unbiased-QCS and $\alpha = 1/(1 + \gamma)$ for MMSE-QCS, and $\boldsymbol{\varepsilon}$ is the scaled quantization noise of $\mathbf{T} \mathbf{z}_t$ (see, e.g., (8)). Therefore,

$$\begin{aligned} \mathbb{E}[\|\mathbf{r}_{t+1}\|_2^2] &= (1 - \beta)^2 \mathbb{E}[\|\mathbf{r}_t\|_2^2] + \mathbb{E}[\|\mathbf{e}_t\|_2^2] \\ &\quad + 2(1 - \beta) \mathbb{E}[\mathbf{r}_t^\top \mathbf{e}_t]. \end{aligned}$$

On the other hand,

$$\begin{aligned} \mathbb{E}[\mathbf{r}_t^\top \mathbf{e}_t] &= \mathbb{E}[\mathbf{r}_t^\top (\mathbf{I} - \alpha \mathbf{T}_t^\top \mathbf{T}_t) \mathbf{z}_t] - \alpha \mathbb{E}[\mathbf{r}_t^\top \varsigma \mathbf{T}_t^\top \boldsymbol{\varepsilon}] \\ &= \mathbb{E}[\mathbb{E}_{\mathbf{T}}[\mathbf{r}_t^\top (\mathbf{I} - \alpha \mathbf{T}_t^\top \mathbf{T}_t) \mathbf{z}_t | \mathbf{r}_t, \dots]] \\ &\quad - \alpha \mathbb{E}[\mathbb{E}_{\boldsymbol{\varepsilon}}[\mathbf{r}_t^\top \varsigma \mathbf{T}_t^\top \boldsymbol{\varepsilon} | \mathbf{T}_t, \mathbf{r}_t, \dots]] \\ &\stackrel{(a)}{=} \mathbb{E}[\mathbb{E}_{\mathbf{T}}[\mathbf{r}_t^\top (\mathbf{I} - \alpha \mathbf{T}_t^\top \mathbf{T}_t) \mathbf{z}_t | \mathbf{r}_t, \dots]] \\ &\stackrel{(b)}{=} (1 - \alpha) \mathbb{E}[\mathbf{r}_t^\top \mathbf{z}_t], \end{aligned}$$

where (a) is because of $\boldsymbol{\varepsilon}$ being an independent zero-mean random vector and (b) due to $\mathbb{E}_{\mathbf{T}}[\mathbf{T}^\top \mathbf{T}] = \mathbf{I}$.

Therefore,

$$\begin{aligned} \mathbb{E}[\|\mathbf{r}_{t+1}\|_2^2] &= (1 - \beta)^2 \mathbb{E}[\|\mathbf{r}_t\|_2^2] + \mathbb{E}[\|\mathbf{e}_t\|_2^2] \\ &\quad + 2(1 - \beta)(1 - \alpha) \mathbb{E}[\mathbf{r}_t^\top \mathbf{z}_t]. \end{aligned}$$

First, we consider the Unbiased-QCS.

Lemma IV. *In Unbiased-QCS with error-feedback, residual signal and the stochastic gradients are uncorrelated, i.e., $\forall t, \tau: \mathbb{E}[\mathbf{g}_t^\top \mathbf{r}_\tau] = 0$.*

Proof. The proof is based on induction. For $\tau = 0$, since $\mathbf{r}_\tau = \mathbf{0}$, the claim holds. Assume that the claim is true for $\tau - 1$.

$$\begin{aligned} \mathbb{E}[\mathbf{g}_t^\top \mathbf{r}_\tau] &= \mathbb{E}[\mathbf{g}_t^\top ((1 - \beta) \mathbf{r}_{\tau-1} + \mathbf{e}_\tau)] \\ &= (1 - \beta) \mathbb{E}[\mathbf{g}_t^\top \mathbf{r}_{\tau-1}] + \mathbb{E}[\mathbf{g}_t^\top \mathbf{e}_\tau] \\ &= \mathbb{E}[\mathbf{g}_t^\top ((\mathbf{I} - \mathbf{T}_\tau^\top \mathbf{T}_\tau) \mathbf{z}_\tau + \varsigma \mathbf{T}_\tau^\top \boldsymbol{\varepsilon}_\tau)] \\ &= \mathbb{E}[\mathbb{E}_{\mathbf{T}_\tau}[\mathbf{g}_t^\top (\mathbf{I} - \mathbf{T}_\tau^\top \mathbf{T}_\tau) \mathbf{z}_\tau | \mathbf{g}_t, \mathbf{z}_\tau]] + \\ &\quad \mathbb{E}[\mathbb{E}_{\boldsymbol{\varepsilon}_\tau}[\varsigma \mathbf{g}_t^\top \mathbf{T}_\tau^\top \boldsymbol{\varepsilon}_\tau | \mathbf{g}_t, \mathbf{T}_\tau]] = 0. \end{aligned}$$

This completes the proof. ■

Since in unbiased QCS, $\alpha = 1$, we have

$$\begin{aligned} \mathbb{E}[\|\mathbf{r}_{t+1}\|_2^2] &= (1 - \beta)^2 \mathbb{E}[\|\mathbf{r}_t\|_2^2] + \mathbb{E}[\|\mathbf{e}_t\|_2^2] \\ &\stackrel{(c)}{\leq} (1 - \beta)^2 \mathbb{E}[\|\mathbf{r}_t\|_2^2] + \gamma \mathbb{E}[\|\mathbf{z}_t\|_2^2] \\ &= (1 - \beta)^2 \mathbb{E}[\|\mathbf{r}_t\|_2^2] + \gamma (\mathbb{E}[\|\mathbf{g}_t\|_2^2] + \beta \mathbb{E}[\|\mathbf{r}_t\|_2^2]) \\ &\stackrel{(d)}{=} ((1 - \beta)^2 + \beta^2 \gamma) \mathbb{E}[\|\mathbf{r}_t\|_2^2] + \gamma \mathbb{E}[\|\mathbf{g}_t\|_2^2] \\ &\stackrel{(e)}{\leq} ((1 - \beta)^2 + \beta^2 \gamma) \mathbb{E}[\|\mathbf{r}_t\|_2^2] + \gamma B \quad (16) \end{aligned}$$

where (c) is due to the fact that for all \mathbf{z} , $\mathbb{E}[\|\mathbf{z} - \hat{\mathbf{z}}\|_2^2] \leq \gamma \|\mathbf{z}\|_2^2$, (d) is from Lemma IV and (e) is from boundedness of \mathbf{g} . The recursive equation (16) with $\mathbf{r}_0 = \mathbf{0}$ implies that

$$\mathbb{E}[\|\mathbf{r}_t\|_2^2] \leq \frac{\gamma}{1 - ((1 - \beta)^2 + \beta^2 \gamma)} B, \quad (17)$$

for all β that $\beta < 1$ and $(1 - \beta)^2 + \beta^2\gamma < 1$, hence, $\beta < \min(1, 2/(1 + \gamma))$.

For MMSE-QCS, the stochastic gradients and residual signal might be correlated. However, for an arbitrary $c > 0$, their correlation can be bounded as

$$\begin{aligned} 0 &\leq \mathbb{E} \left[\left\| \frac{1}{\sqrt{c}} \mathbf{g}_t \pm \sqrt{c} \mathbf{r}_\tau \right\|_2^2 \right] = \frac{1}{c} \mathbb{E} [\|\mathbf{g}_t\|_2^2] + c \mathbb{E} [\|\mathbf{r}_\tau\|_2^2] \\ &\quad \pm 2 \mathbb{E} [\mathbf{g}_t^\top \mathbf{r}_\tau] \\ &\Rightarrow 2 |\mathbb{E} [\mathbf{g}_t^\top \mathbf{r}_\tau]| \leq \frac{1}{c} \mathbb{E} [\|\mathbf{g}_t\|_2^2] + c \mathbb{E} [\|\mathbf{r}_\tau\|_2^2]. \end{aligned}$$

Therefore, noting that $\mathbb{E} [\|\mathbf{e}\|_2^2] \leq (1 - \alpha) \|\mathbf{z}\|_2^2$ and $\mathbb{E} [\|\mathbf{g}_t\|_2^2] \leq B$,

$$\begin{aligned} \mathbb{E} [\|\mathbf{r}_{t+1}\|_2^2] &= (1 - \beta)^2 \mathbb{E} [\|\mathbf{r}_t\|_2^2] + \mathbb{E} [\|\mathbf{e}_t\|_2^2] \\ &\quad + 2(1 - \beta)(1 - \alpha) \mathbb{E} [\mathbf{r}_t^\top \mathbf{z}_t] \\ &\leq (1 - \beta)^2 \mathbb{E} [\|\mathbf{r}_t\|_2^2] + (1 - \alpha) \mathbb{E} [\|\mathbf{g}_t + \beta \mathbf{r}_t\|_2^2] + \\ &\quad 2(1 - \beta)(1 - \alpha) \mathbb{E} [\mathbf{r}_t^\top (\mathbf{g}_t + \beta \mathbf{r}_t)] \\ &\leq (1 - \alpha\beta(2 - \beta)) \mathbb{E} [\|\mathbf{r}_t\|_2^2] + \\ &\quad (1 - \alpha)(c \mathbb{E} [\|\mathbf{r}_t\|_2^2] + \frac{1}{c} B) + (1 - \alpha) B \\ &= (1 - \alpha\beta(2 - \beta) + c(1 - \alpha)) \mathbb{E} [\|\mathbf{r}_t\|_2^2] + \\ &\quad (1 - \alpha)(1 + \frac{1}{c}) B. \end{aligned}$$

Therefore, if $|1 - \alpha\beta(2 - \beta) + c(1 - \alpha)| < 1$,

$$\mathbb{E} [\|\mathbf{r}_t\|_2^2] \leq \frac{1 + 1/c}{\alpha\beta(2 - \beta) - c(1 - \alpha)} (1 - \alpha) B. \quad (18)$$

Minimizing w.r.t. c and substituting $\alpha = 1/(1 + \gamma)$ results in

$$\mathbb{E} [\|\mathbf{r}_t\|_2^2] \leq \frac{1 - \alpha}{(\sqrt{1 - \alpha(1 - \beta)^2} - \sqrt{1 - \alpha})^2} B \quad (19)$$

$$= \frac{\gamma}{(\sqrt{\gamma + 1 - (1 - \beta)^2} - \sqrt{\gamma})^2} B. \quad (20)$$

Note that if the SGs have bounded variance, i.e., $\mathbb{E} [\|\mathbf{g} - \nabla f\|_2^2] \leq \sigma^2$, a similar approach can be used to bound $\mathbb{E} [\|\mathbf{r}_t\|_2^2]$ based on the σ^2 and the weighted average of $\|\nabla f(\mathbf{w}_{t-i})\|_2^2$ for $i = 0, \dots, t$. This is specially helpful when analyzing the convergence of the training algorithm with error feedback under the assumption of bounded variance SG.

6 Convergence Analysis

Proof of Lemm 5

The proof follows the same line of argument as for ordinary SGD which is repeated here for the sake of completeness.

Recall that for Lipschitz-smooth function $f(\cdot)$, for arbitrary \mathbf{w} and δ ,

$$f(\mathbf{w} + \delta) \leq f(\mathbf{w}) + \delta^\top \nabla f(\mathbf{w}) + \frac{L}{2} \|\delta\|_2^2.$$

First, we consider Unbiased-QCS. Tt the t -th iteration of training, $\mathbf{w}_{t+1} = \mathbf{w}_t - \mu \hat{\mathbf{g}}_t$, where $\mathbb{E} [\hat{\mathbf{g}}_t] = \nabla f(\mathbf{w}_t)$ and $\mathbb{E} [\|\hat{\mathbf{g}}_t\|_2^2] \leq (1 + \gamma) \mathbb{E} [\|\mathbf{g}\|_2^2] \leq (1 + \gamma) B$. Hence,

$$\begin{aligned} \mathbb{E} [f(\mathbf{w}_{t+1})] &\leq f(\mathbf{w}_t) + \langle \nabla f(\mathbf{w}_t), \mathbb{E} [\mathbf{w}_{t+1} - \mathbf{w}_t] \rangle \\ &\quad + \frac{L}{2} \mathbb{E} [\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2] \\ &\leq f(\mathbf{w}_t) - \mu \|\nabla f(\mathbf{w}_t)\|_2^2 + \frac{L}{2} \mu^2 (1 + \gamma) B. \end{aligned}$$

Rearranging terms, taking expectation and summing from $t = 0$ to $T - 1$, results in

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\|\nabla f(\mathbf{w}_t)\|_2^2] &\leq \frac{f(\mathbf{w}_0) - \mathbb{E} [f(\mathbf{w}_T)]}{T\mu} + \frac{L}{2} \mu (1 + \gamma) B \\ &\leq \frac{f(\mathbf{w}_0) - f^*}{T\mu} + \frac{L}{2} \mu (1 + \gamma) B. \end{aligned}$$

Setting $\mu = 1/\sqrt{T}$, results in

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\|\nabla f(\mathbf{w}_t)\|_2^2] \leq \frac{f(\mathbf{w}_0) - f^* + \frac{L}{2} (1 + \gamma) B}{\sqrt{T}}.$$

Note that in the case that the stochastic gradients have bounded variance³, i.e., $\mathbb{E} [\|\mathbf{g} - \nabla f\|_2^2] \leq \sigma^2$, $\mathbb{E} [\|\hat{\mathbf{g}}_t\|_2^2] \leq (1 + \gamma) \mathbb{E} [\|\mathbf{g}\|_2^2] \leq (1 + \gamma)(\sigma^2 + \|\nabla f\|_2^2)$ and we can modify the above argument as follows to bound the convergence rate,

$$\begin{aligned} \mathbb{E} [f(\mathbf{w}_{t+1})] - f(\mathbf{w}_t) &\leq + \langle \nabla f(\mathbf{w}_t), \mathbb{E} [\mathbf{w}_{t+1} - \mathbf{w}_t] \rangle \\ &\quad + \frac{L}{2} \mathbb{E} [\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2] \\ &\leq -\mu \|\nabla f(\mathbf{w}_t)\|_2^2 + \frac{L}{2} \mu^2 (1 + \gamma) (\sigma^2 + \|\nabla f(\mathbf{w}_t)\|_2^2) \\ &= -(\mu - \frac{L}{2} \mu^2 (1 + \gamma)) \|\nabla f(\mathbf{w}_t)\|_2^2 + \frac{L}{2} \mu^2 (1 + \gamma) \sigma^2. \end{aligned}$$

Following same argument as before,

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\|\nabla f(\mathbf{w}_t)\|_2^2] &\leq \frac{2(f(\mathbf{w}_0) - f^*)}{T(2\mu - L\mu^2(1 + \gamma))} \\ &\quad + \frac{L\mu^2(1 + \gamma)}{2\mu - L\mu^2(1 + \gamma)} \sigma^2. \end{aligned}$$

It can be verified that if $T > 4L^2(\gamma + 1)^2$, we can find $\mu \leq 2/\sqrt{T}$ such that $2\mu - L\mu^2(\gamma + 1) = 2/\sqrt{T}$. This simplifies the above equation to

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\|\nabla f(\mathbf{w}_t)\|_2^2] \leq \frac{f(\mathbf{w}_0) - f^*}{\sqrt{T}} + \frac{2L(1 + \gamma)}{\sqrt{T}} \sigma^2.$$

The analysis for MMSE-QCS is straightforward. Note that $\hat{\mathbf{g}}_{mmse} = \frac{1}{\gamma+1} \hat{\mathbf{g}}_u$, where $\hat{\mathbf{g}}_{mmse}$ is the MMSE-QCS quantized SG and $\hat{\mathbf{g}}_u$ is the output of Unbiased-QCS. Therefore, training with MMSE-QCS and step-size μ would be the same as using Unbiased-QCS with step-size $\mu/(\gamma + 1)$.

³In this case, it is not necessary to assume that the cost function has bounded gradient everywhere.

Remark 3. Note that since Unbiased-QCS has bounded variance and is unbiased, the compressed SG will be stochastic gradient itself with bounded variance. Hence, majority of the results can be readily applied to prove the convergence of Unbiased-QCS and MMSE-QCS under different conditions such as (Bottou 1998; Bubeck 2015).

Proof of Lemma 6

The proof is based on the ideas from (Karimireddy et al. 2019) and follows the similar arguments with slight modifications, which is repeated here for the sake of completeness.

Let $\tilde{\mathbf{w}}_t = \mathbf{w}_t - \mu \mathbf{r}_t$. Note that since by Lemma ?? the residue signal has bounded variance, $\tilde{\mathbf{w}}_t$ would be bounded. It can be easily verified that $\tilde{\mathbf{w}}_{t+1} = \tilde{\mathbf{w}}_t - \mu \mathbf{g}_t$. Hence, following similar argument as in (Karimireddy et al. 2019), for arbitrary $\rho > 0$

$$\begin{aligned} & \mathbb{E}[f(\tilde{\mathbf{w}}_{t+1})] - f(\tilde{\mathbf{w}}_t) \\ & \leq \frac{L}{2} \mathbb{E}[\|\tilde{\mathbf{w}}_{t+1} - \tilde{\mathbf{w}}_t\|_2^2] + \langle \nabla f(\tilde{\mathbf{w}}_t), \mathbb{E}[\tilde{\mathbf{w}}_{t+1} - \tilde{\mathbf{w}}_t] \rangle \\ & \leq \frac{L}{2} \mu^2 B - \mu(1 - \rho) \|\nabla f(\mathbf{w}_t)\|_2^2 + \frac{1}{\rho} L^2 \mu^3 \mathbb{E}[\|\mathbf{r}_t\|_2^2]. \end{aligned}$$

On the other hand, from Lemma ??, the residue is bounded as $\mathbb{E}[\|\mathbf{r}_t\|_2^2] \leq \eta B$ where η is a constant depending on the β (weight of error feedback) and γ , according to (??) or (??) of the main paper for Unbiased-QCS and MMSE-QCS. Therefore,

$$\begin{aligned} \mu(1 - \rho) \|\nabla f(\mathbf{w}_t)\|_2^2 & \leq \left(\frac{L}{2} \mu^2 + \frac{1}{\rho} L^2 \mu^3 \eta \right) B \\ & \quad + f(\tilde{\mathbf{w}}_t) - \mathbb{E}[f(\tilde{\mathbf{w}}_{t+1})]. \end{aligned}$$

Taking expectation, rearranging terms and noting that $\tilde{\mathbf{w}}_0 = \mathbf{w}_0$ and $\mathbb{E}[f(\mathbf{w}_T)] \geq f^*$, we conclude that for $0 < \rho < 1$,

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(\mathbf{w}_t)\|_2^2] \leq \frac{f(\mathbf{w}_0) - f^*}{T\mu(1 - \rho)} + \frac{LB}{1 - \rho} \left(\frac{\mu}{2} + \frac{L\eta\mu^2}{\rho} \right).$$

Setting $\rho = 0.5$, gives the desired result.

For tighter analysis, let μ and ρ be such that $\mu(1 - \rho) = 1/\sqrt{T}$ and $\frac{\mu}{1 - \rho} = \frac{1 + \epsilon}{\sqrt{T}}$ for arbitrary $\epsilon > 0$, i.e.,

$$\mu = \frac{\sqrt{1 + \epsilon}}{\sqrt{T}}, \quad 1 - \rho = \frac{1}{\sqrt{1 + \epsilon}}.$$

Therefore,

$$\begin{aligned} \min_t \mathbb{E}[\|\nabla f(\mathbf{w}_t)\|_2^2] & \leq \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(\mathbf{w}_t)\|_2^2] \\ & \leq \frac{f(\mathbf{w}_0) - f^* + \frac{L}{2} B}{\sqrt{T}} + L^2 B \frac{(1 + \epsilon)^2}{\sqrt{1 + \epsilon} - 1} \frac{\eta}{T}. \end{aligned}$$

References

Alistarh, D.; Grubic, D.; Li, J.; Tomioka, R.; and Vojnovic, M. 2017. QSGD: Communication-Efficient SGD via Gradient Quantization and Encoding. In *Advances in Neural Information Processing Systems*, 1707–1718.

Bottou, L. 1998. Online Learning and Stochastic Approximations, Revised 2018. *On-Line Learning in Neural Networks* 17(9):1–35.

Bourgain, J.; Vu, V. H.; and Wood, P. M. 2010. On the singularity probability of discrete random matrices. *Journal of Functional Analysis* 258(2):559–603.

Bubeck, S. 2015. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning* 8(3–4):231–357.

Gray, R. M., and Neuhoff, D. L. 1998. Quantization. *IEEE Transactions on Information Theory* 44(6):2325–2383.

Gray, R. M., and Stockham, T. G. 1993. Dithered quantizers. *IEEE Transactions on Information Theory* 39(3):805–812.

Karimireddy, S. P.; Rebjock, Q.; Stich, S. U.; and Jaggi, M. 2019. Error Feedback Fixes SignSGD and other Gradient Compression Schemes. *arXiv preprint arXiv:1901.09847*.

Schuchman, L. 1964. Dither signals and their effect on quantization noise. *IEEE Transactions on Communication Technology* 12(4):162–165.

Wangni, J.; Wang, J.; Liu, J.; and Zhang, T. 2018. Gradient sparsification for communication-efficient distributed optimization. In *Advances in Neural Information Processing Systems*, 1306–1316.

Wen, W.; Xu, C.; Yan, F.; Wu, C.; Wang, Y.; Chen, Y.; and Li, H. 2017. TernGrad: Ternary Gradients to Reduce Communication in Distributed Deep Learning. In *TernGrad: Ternary Gradients to Reduce Communication in Distributed Deep Learning*, 1509–1519. Curran Associates, Inc.