

---

# Quantized Stochastic Gradient Descent: Communication versus Convergence

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Parallel implementations of stochastic gradient descent (SGD) have received signif-  
2 icant research attention, thanks to excellent scalability properties of this algorithm,  
3 and to its efficiency in the context of training deep neural networks. A fundamental  
4 barrier for parallelizing large-scale SGD is the fact that the cost of communicat-  
5 ing the gradient updates between nodes can be very large. Consequently, lossy  
6 compression heuristics have been proposed, by which nodes only communicate  
7 quantized gradients. Although effective in practice, these heuristics do not always  
8 provably converge, and it is not clear whether they are optimal. In this paper, we  
9 propose *Quantized SGD (QSGD)*, a family of compression schemes which allow  
10 the compression of gradient updates at each node, while guaranteeing convergence  
11 under standard assumptions. QSGD allows the user to trade off compression and  
12 convergence time: it can communicate a *sublinear* number of bits per iteration  
13 in the model dimension, and can achieve asymptotically optimal communication  
14 cost. We complement our theoretical results with empirical data, showing that  
15 QSGD can significantly reduce communication cost, while being competitive with  
16 standard uncompressed techniques on a variety of real tasks.

## 17 1 Introduction

18 Let  $\mathcal{X} \subseteq \mathbb{R}^n$  be a known convex set, and consider stochastic gradient descent for a smooth function  
19  $f : \mathcal{X} \rightarrow \mathbb{R}$ , in which we only have access to independent stochastic gradients of  $f$ . We assume that  
20 stochastic gradient  $\tilde{g}(\mathbf{x})$  is unbiased  $\mathbb{E}[\tilde{g}(\mathbf{x})] = \nabla f(\mathbf{x})$  and satisfies the second moment condition  
21  $\mathbb{E}[\|\tilde{g}(\mathbf{x})\|_2^2] \leq B$  for all  $\mathbf{x} \in \mathcal{X}$ .

22 Now consider a synchronous parallel stochastic gradient descent setting in which we have  $K$  workers,  
23 each of which have access to independent stochastic gradients of  $f$ . Each worker computes the  
24 stochastic gradient synchronously and communicates the gradients with each other. After the  
25 communication, each worker updates the parameter using the aggregated gradient as

$$\mathbf{x}_{t+1} = \Pi_{\mathcal{X}} \left( \mathbf{x}_t - \frac{\eta_t}{K} \sum_{\ell=1}^K \tilde{g}^{\ell}(\mathbf{x}_t) \right),$$

26 where  $\tilde{g}^{\ell}(\mathbf{x}_t)$  is the stochastic gradient computed on the  $\ell$ th worker.

27 One can easily imagine that when the number of parameters  $n$  is large, the cost of communication  
28 can be significant. 1-Bit SGD [5] addresses this issue by introducing a quantization function that  
29 roughly speaking encodes a gradient vector into one bit for each coordinate corresponding to its sign  
30 and two float numbers corresponding to the mean of the positive coordinates and the mean of the

31 negative coordinates. On the receiver’s side the gradient vector can be approximately recovered by  
 32 setting all the positive coordinates to the positive mean and the negative coordinates as the negative  
 33 mean. This heuristic requires roughly  $n$  bits per gradient, but does not guarantee convergence.

## 34 2 A Randomized Quantization Scheme

35 We propose the random quantization function  $Q(\mathbf{v})$  defined as follows:

$$Q_i(\mathbf{v}) = \|\mathbf{v}\|_2 \cdot \text{sgn}(v_i) \xi_i(\mathbf{v}), \quad (1)$$

36 where  $\xi_i(\mathbf{v})$ ’s are independent random variables such that  $\xi_i(\mathbf{v}) = 1$  with probability  $|v_i|/\|\mathbf{v}\|_2$ , and  
 37  $\xi_i(\mathbf{v}) = 0$ , otherwise. If  $\mathbf{v} = \mathbf{0}$ , we define  $Q(\mathbf{v}) = \mathbf{0}$ .

38 The key properties of  $Q[\tilde{g}(\mathbf{x})]$  are sparsity, unbiasedness, and bounded second moment as shown in  
 39 the following lemma:

40 **Lemma 2.1.** *For any  $\mathbf{v} \in \mathbb{R}^n$ , we have  $\mathbb{E}[\|Q(\mathbf{v})\|_0] \leq \sqrt{n}$  (sparsity),  $\mathbb{E}[Q(\mathbf{v})] = \mathbf{v}$  (unbiasedness),  
 41 and  $\mathbb{E}[\|Q(\mathbf{v})\|^2] \leq \sqrt{n}\|\mathbf{v}\|_2^2$  (second moment bound).*

42 The sparsity allows us to succinctly encode  $Q(\mathbf{x})$ , for any  $\mathbf{x}$ , in expectation. The information  
 43 contained in  $Q(\mathbf{v})$  can be expressed by (1) a float variable that encodes the value of  $\|\mathbf{v}\|_2$ , (2)  
 44 identities of the vector coordinates  $i$  for which  $\xi_i(\mathbf{v}) = 1$ , and (3) the values of signs  $\text{sgn}(v_i)$  for  
 45 these coordinates. Let  $\text{Code}(Q(\mathbf{v}))$  denote a binary representation of such a tuple representation of  
 46  $Q(\mathbf{v})$ . Then, one can show the following bound, whose proof is deferred to the full version of our  
 47 paper.

48 **Lemma 2.2.** *For every vector  $\mathbf{v} \in \mathbb{R}^n$ , we have  $\mathbb{E}[\|\text{Code}(Q(\mathbf{v}))\|] \leq \sqrt{n}(\log(n) + \log(2e)) + F$ ,  
 49 where  $F$  is the number of bits for representing one floating point number.*

50 These two lemmas together imply the following theorem.

51 **Theorem 2.3.** *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be fixed, and let  $\mathbf{x} \in \mathbb{R}^n$  be arbitrary. If  $\tilde{g}(\mathbf{x})$  is a stochastic  
 52 gradient for  $f$  at  $\mathbf{x}$  with second moment bound  $B$ , then  $Q(\tilde{g}(\mathbf{x}))$  is a stochastic gradient for  $f$  at  $\mathbf{x}$   
 53 with second moment bound  $\sqrt{n}B$ . Moreover, in expectation  $Q(\tilde{g}(\mathbf{x}))$  can be communicated using  
 54  $\sqrt{n}(\log n + \log 2e) + F$  bits.*

55 Note that the above communication cost is *sublinear* in the dimension  $n$  compared to the linear cost  
 56 when the gradients are communicated without compression or using 1-bit SGD. Using standard con-  
 57 vergence results (e.g., [1]), we obtain a stochastic gradient algorithm that requires only  $O(\sqrt{n} \log n)$   
 58 communication per round and converges to the same precision in  $O(\sqrt{n})$  times more iterations.

59 We can control the trade-off between communication and convergence by introducing *bucketing*.  
 60 More precisely, we partition the gradient vector into  $n/d$  buckets, each of which containing  $d$   
 61 consecutive coordinates, and apply the quantization and encoding to each bucket. Then a simple  
 62 extension of Theorem 2.3 predicts that the second moment bound then becomes  $\sqrt{d}B$ . Setting  $d = 1$ ,  
 63 we recover no quantization (vanilla SGD), and  $d = n$  corresponds to full quantization. However,  
 64 since quantization increases the second moment bound, there is no improvement in terms of the total  
 65 communication cost.

## 66 3 A Generalized Randomized Quantization Scheme

67 In order to explore the trade-off between communication and convergence more carefully, we propose  
 68 a more general lossy-compression scheme defined as follows:

$$Q_i(\mathbf{v}, s) = \|\mathbf{v}\|_2 \cdot \text{sgn}(v_i) \xi_i(\mathbf{v}, s), \quad (2)$$

69 where  $s \geq 1$  is a tuning parameter,  $\xi_i(\mathbf{v}, s)$ ’s are independent random variables with distributions  
 70 defined as follows. Let  $0 \leq \ell < s$  be an integer such that  $|v_i|/\|\mathbf{v}\|_2 \in [\ell/s, (\ell + 1)/s]$ . Then

$$\xi_i(\mathbf{v}, s) = \begin{cases} \ell/s & \text{with probability } 1 - p\left(\frac{|v_i|}{\|\mathbf{v}\|_2}, s\right); \\ (\ell + 1)/s & \text{otherwise.} \end{cases}$$

71 Here,  $p(a, s) = as - \ell$  for any  $a \in [0, 1]$ . If  $\mathbf{v} = \mathbf{0}$ , then we define  $Q(\mathbf{v}, s) = \mathbf{0}$ .

72 The random quantization function (1) corresponds to the special case  $s = 1$ . We obtain the three key  
73 properties as we show in the following lemma.

74 **Lemma 3.1.** *For any  $\mathbf{v} \in \mathbb{R}^n$ , we have that  $\mathbb{E}[\|Q(\mathbf{v}, s)\|_0] \leq s^2 + \sqrt{n}$  (sparsity),  $\mathbb{E}[Q(\mathbf{v}, s)] = \mathbf{v}$   
75 (unbiasedness), and  $\mathbb{E}[\|Q(\mathbf{v}, s)\|_2^2] \leq (1 + \min(n/s^2, \sqrt{n}/s))\|\mathbf{v}\|_2^2$  (second moment bound).*

76 Note that the factor  $c_{n,s} := 1 + \min(n/s^2, \sqrt{n}/s)$  in the second-moment bound is parameterized  
77 with the dimension  $n$  and the tuning parameter  $s$ . For the special case  $s = 1$ , we have  $c_{n,s} = \Theta(\sqrt{n})$ ,  
78 which is consistent with the result in Lemma 2.1. By varying the value of the parameter  $s$  between 1  
79 and  $\sqrt{n}$ , we can smoothly vary  $c_{n,s}$  between  $\Theta(\sqrt{n})$  and  $\Theta(1)$ . We can also see that as we increase  
80  $s$ , the quantized gradient becomes less sparse. The sparsity bound is  $O(\sqrt{n})$  at  $s = 1$  and  $O(n)$   
81 at  $s = \sqrt{n}$ . We also note that the distribution of  $\xi_i(\mathbf{v}, s)$  is a unique distribution that has minimal  
82 variance over distributions that have support  $\{0, 1/s, \dots, 1\}$  and unbiased.

83 In the sparse regime where we expect the quantized gradient to contain at most  $n/2$  non-zero  
84 coordinates, we have the following theorem.

85 **Theorem 3.2.** *Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be fixed, and let  $\mathbf{x} \in \mathbb{R}^n$  be arbitrary. If  $\tilde{g}(\mathbf{x})$  is a stochastic gradient  
86 for  $f$  at  $\mathbf{x}$  with second moment bound  $B$ , then  $Q_s(\tilde{g}(\mathbf{x}))$  is a stochastic gradient for  $f$  at  $\mathbf{x}$  with  
87 second moment bound  $(1 + \min(n/s^2, \sqrt{n}/s))B$ . Moreover, there is an encoding scheme so that  
88 in expectation, the number of bits needed to communicate  $Q_s(\tilde{g}(\mathbf{x}))$  is upper bounded by*

$$F + \left(3 + \frac{3}{2} \cdot (1 + o(1)) \log \left( \frac{2(s^2 + n)}{s^2 + \sqrt{n}} \right) \right) (s^2 + \sqrt{n}).$$

89 The communication cost can be, roughly speaking, broken down into one float number representing  
90 the norm and  $s^2 + \sqrt{n}$  (in expectation) bits and integers representing the signs, magnitudes, and  
91 positions of the non-zero coordinates. We use the *recursive Elias coding*, which is favorable for small  
92 integers, to achieve the above bound.

93 For large  $s$ , the quantized gradient becomes dense, and we no longer need to communicate the  
94 positions of the non-zero coordinates.

95 **Theorem 3.3.** *Let  $f$ ,  $\mathbf{x}$ , and  $\tilde{g}(\mathbf{x})$  be as in Theorem 3.2. There is an encoding scheme for  $Q_s(\tilde{g}(\mathbf{x}))$   
96 which in expectation has length*

$$F + \left( \frac{1 + o(1)}{2} \left( \log \left( 1 + \frac{s^2 + \min(n, s\sqrt{n})}{n} \right) + 1 \right) + 2 \right) n.$$

97 *In particular, if  $s = \sqrt{n}$ , then this encoding requires  $\leq F + 2.8n$  bits in expectation.*

98 In fact, for  $s = \sqrt{n}$ , the second moment bound is only 2 times worse than no quantization and the  
99 communication cost is only 2.8 bits per coordinate. We defer the description of the quantization  
100 schemes in Theorems 3.2 and 3.3, and their use in the context of SVRG [2], to the full version.

## 101 4 Experiments

102 We now empirically validate our approach, using experiments aimed at data-parallel and model-  
103 parallel settings. We have implemented QSGD on GPUs using deep learning framework Chainer [7].

104 **Quantization vs. Accuracy.** In the first set of experiments, we explore the relation between  
105 performance and the granularity at which quantization is applied to the gradient vector.

106 Here, our experiments deviate from the theory, as we use a deep network, with non-convex objective.

107 *MNIST dataset.* The first dataset is the MNIST dataset of handwritten digits. The training set consists  
108 of 60,000 28 x 28 single digit images. The test set consists of 10,000 images. We train a two-layer  
109 perceptron with 4096 hidden units and ReLU activation with a minibatch size of 256 and step size  
110 of 0.1. Results are shown in Figure 1(a). Rather surprisingly, in terms of both training negative  
111 log-likelihood loss and the test accuracy, QSGD *improves* performance. This is consistent with  
112 recent work [4] suggesting benefits of added noise in training deep networks. We observed no such  
113 improvement for a linear model on the same dataset.

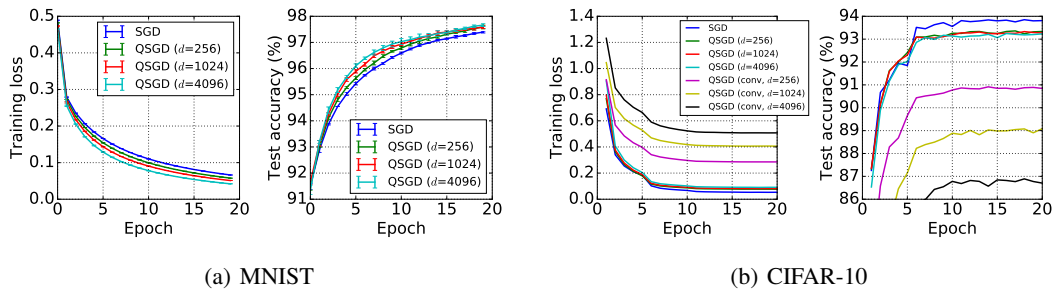


Figure 1: Training on a single machine on MNIST and CIFAR-10. SGD corresponds to bucket size of  $d = 1$ . QSGD performs better in terms of both training loss and test accuracy on the MNIST dataset.

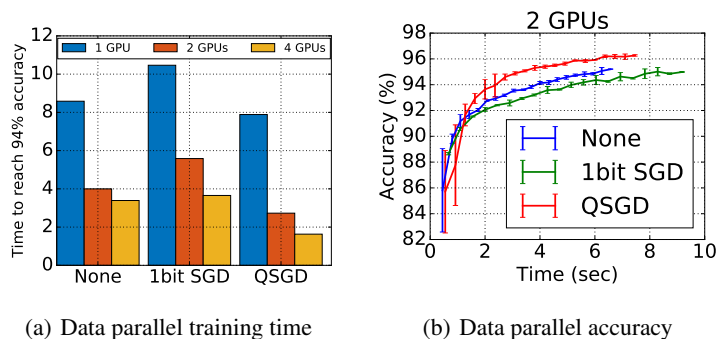


Figure 2: Multi GPU experiment.

114 The total number of parameters of this model is 3.3 million, most of them lying in the first layer.  
 115 Using Theorem 2.2, we can approximate the effective number of floats communicated by QSGD.  
 116 Assuming  $F = 32$ , we get roughly 88k, 49k, and 29k effective floats for bucket sizes  $d = 256, 1024,$   
 117 and 4096, respectively. There is a massive reduction in communication since for each bucket we only  
 118 need to communicate one float and the positions and signs of  $\tilde{O}(\sqrt{d})$  entries, each of which only  
 119 requires  $O(\log d)$  bits, which is typically much smaller than 32 (e.g., 11 bits for  $d = 256$ ).

120 *CIFAR-10 dataset.* Next, we consider the CIFAR-10 object classification dataset [3]. The original  
 121 training set consists of 50,000  $32 \times 32$  color images, augmented by translating, cropping with window  
 122 size  $28 \times 28$ , and horizontal flipping. The augmented training set contains 1.8 million images.

123 We use a small VGG model [6] consisting of nine 2D convolution layers and three fully connected  
 124 layers. The total number of parameters is roughly 22 million. All methods used momentum of 0.9.  
 125 See the full paper for the details. When we only quantized the fully connected layers, we have found  
 126 that the bucket size can be increased without much loss in accuracy (see Fig. 1(b)). The effective  
 127 number of floats to be communicated are 1.5 million, 1.3 million, and 1.2 million for bucket sizes  
 128  $d = 256, 1024,$  and 4096, respectively. On the other hand, when we also applied the quantization  
 129 to the convolutional layers, we observed a noticeable increase in the training objective as well as  
 130 reduction in the test accuracy. The effective number of floats to be communicated are 580k, 312k,  
 131 176k, respectively.

132 **Parallelization.** In Figure 2 (a) and (b), we show preliminary scalability experiments on MNIST,  
 133 using up to 4 GPUs, compared with vanilla SGD and 1-Bit SGD [5]. The setup is the same as in  
 134 the previous section, and we use double buffering [5] to perform communication and quantization  
 135 concurrently with the computation. Experiments are preliminary in the sense that we did not fully  
 136 optimize either 1-Bit SGD or QSGD to their full potential; in particular, quantized gradients are  
 137 communicated in raw floats instead of using more efficient encoding.

## 138 **References**

- 139 [1] S. Bubeck. Convex optimization: Algorithms and complexity. *Foundations and Trends® in*  
140 *Machine Learning*, 8(3-4):231–357, 2015.
- 141 [2] R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance  
142 reduction. In *NIPS*, 2013.
- 143 [3] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images, 2009.
- 144 [4] A. Neelakantan, L. Vilnis, Q. V. Le, I. Sutskever, L. Kaiser, K. Kurach, and J. Martens. Adding  
145 gradient noise improves learning for very deep networks. *arXiv preprint arXiv:1511.06807*,  
146 2015.
- 147 [5] F. Seide, H. Fu, J. Droppo, G. Li, and D. Yu. 1-bit stochastic gradient descent and its application  
148 to data-parallel distributed training of speech dnns. In *INTERSPEECH*, 2014.
- 149 [6] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image  
150 recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- 151 [7] S. Tokui, K. Oono, S. Hido, C. San Mateo, and J. Clayton. Chainer: a next-generation open  
152 source framework for deep learning.