# **Qsparse-local-SGD: Distributed SGD with Quantization, Sparsification, and Local Computations**

**Debraj Basu** \* Adobe Inc. dbasu@adobe.com

Deepesh Data UCLA deepeshdata@ucla.edu

Can Karakus \*
Amazon Inc.
cakarak@amazon.com

Suhas Diggavi UCLA suhasdiggavi@ucla.edu

## **Abstract**

Communication bottleneck has been identified as a significant issue in distributed optimization of large-scale learning models. Recently, several approaches to mitigate this problem have been proposed, including different forms of gradient compression or computing local models and mixing them iteratively. In this paper we propose *Qsparse-local-SGD* algorithm, which combines aggressive sparsification with quantization and local computation along with error compensation, by keeping track of the difference between the true and compressed gradients. We propose both synchronous and asynchronous implementations of *Qsparse-local-SGD*. We analyze convergence for *Qsparse-local-SGD* in the *distributed* case, for smooth non-convex and convex objective functions. We demonstrate that *Qsparse-local-SGD* converges at the same rate as vanilla distributed SGD for many important classes of sparsifiers and quantizers. We use *Qsparse-local-SGD* to train ResNet-50 on ImageNet, and show that it results in significant savings over the state-of-the-art, in the number of bits transmitted to reach target accuracy.

## 1 Introduction

Stochastic Gradient Descent (SGD) [14] and its many variants have become the workhorse for modern large-scale optimization as applied to machine learning [5,8]. We consider the setup where SGD is applied to the *distributed* setting, where R different nodes compute *local* SGD on their *own* datasets  $\mathcal{D}_r$ . Co-ordination between them is done by aggregating these local computations to update the overall parameter  $\mathbf{x}_t$  as,  $\mathbf{x}_{t+1} = \mathbf{x}_t - \frac{\eta_t}{R} \sum_{r=1}^R g_t^r$ , where  $\{g_t^r\}_{r=1}^R$  are the local stochastic gradients at the R machines for a local loss function  $f^{(r)}(\mathbf{x})$  of the parameters, where  $f^{(r)}: \mathbb{R}^d \to \mathbb{R}$ .

It is well understood by now that sending full-precision gradients, causes communication to be the bottleneck for many large scale models [4,7,33,39]. The communication bottleneck could be significant in emerging edge computation architectures suggested by federated learning [1,17,22]. To address this, many methods have been proposed recently, and these methods are broadly based on three major approaches: (i) *Quantization* of gradients, where nodes locally quantize the gradient (perhaps with randomization) to a small number of bits [3,7,33,39,40]. (ii) *Sparsification* of gradients, *e.g.*, where nodes locally select  $Top_k$  values of the gradient in absolute value and transmit these at full precision [2,4,20,30,32,40], while maintaining errors in local nodes for later compensation. (iii) *Skipping communication rounds* whereby nodes average their models after locally updating their models for several steps [9,10,31,34,37,43,45].

<sup>\*</sup>Work done while Debraj Basu and Can Karakus were at UCLA.

In this paper we propose *Qsparse-local-SGD* algorithm, which combines aggressive sparsification with quantization and local computation along with error compensation, by keeping track of the difference between the true and compressed gradients. We propose both synchronous and asynchronous<sup>2</sup> implementations of *Qsparse-local-SGD*. We analyze convergence for *Qsparse-local-SGD* in the *distributed* case, for smooth non-convex and convex objective functions. We demonstrate that, *Qsparse-local-SGD* converges at the same rate as vanilla distributed SGD for many important classes of sparsifiers and quantizers. We implement *Qsparse-local-SGD* for ResNet-50 using the ImageNet dataset, and show that we achieve target accuracies with a small penalty in final accuracy (approximately 1 %), with about a factor of 15-20 savings over the state-of-the-art [4, 30, 31], in the total number of bits transmitted. While the downlink communication is not our focus in this paper (also in [4, 20, 39], for example), it can be inexpensive when the broadcast routine is implemented in a tree-structured manner as in many MPI implementations, or if the parameter server aggregates the sparse quantized updates and broadcasts it.

Related work. The use of quantization for communication efficient gradient methods has decades rich history [11] and its recent use in training deep neural networks [27, 32] has re-ignited interest. Theoretically justified gradient compression using unbiased stochastic quantizers has been proposed and analyzed in [3, 33, 39]. Though methods in [36, 38] use induced sparsity in the quantized gradients, explicitly sparsifying the gradients more aggressively by retaining  $Top_k$  components, e.g., k < 1%, has been proposed [2, 4, 20, 30, 32], combined with error compensation to ensure that all co-ordinates do get eventually updated as needed. [40] analyzed error compensation for QSGD, without  $Top_k$  sparsification and a focus on quadratic functions. Another approach for mitigating the communication bottlenecks is by having infrequent communication, which has been popularly referred to in the literature as iterative parameter mixing and model averaging, see [31,43] and references therein. Our work is most closely related to and builds on the recent theoretical results in [4, 30, 31, 43]. [30] considered the analysis for the centralized  $Top_k$  (among other sparsifiers), and [4] analyzed a distributed version with the assumption of closeness of the aggregated  $Top_k$ gradients to the centralized  $Top_k$  case, see Assumption 1 in [4], [31,43] studied local-SGD, where several local iterations are done before sending the full gradients, and did not do any gradient compression beyond local iterations. Our work generalizes these works in several ways. We prove convergence for the distributed sparsification and error compensation algorithm, without the assumption of [4], by using the perturbed iterate methods [21, 30]. We analyze non-convex (smooth) objectives as well as strongly convex objectives for the distributed case with local computations. [30] gave a proof only for convex objective functions and for centralized case and therefore without local computations<sup>3</sup>. Our techniques compose a (stochastic or deterministic 1-bit sign) quantizer with sparsification and local computations using error compensation; in fact this technique works for any compression operator satisfying a regularity condition (see Definition 3).

Contributions. We study a distributed set of R worker nodes each of which perform computations on locally stored data denoted by  $\mathcal{D}_r$ . Consider the empirical-risk minimization of the loss function  $f(\mathbf{x}) = \frac{1}{R} \sum_{r=1}^R f^{(r)}(\mathbf{x})$ , where  $f^{(r)}(\mathbf{x}) = \underset{i \sim \mathcal{D}_r}{\mathbb{E}} [f_i(\mathbf{x})]$ , where  $\underset{i \sim \mathcal{D}_r}{\mathbb{E}} [\cdot]$  denotes expectation<sup>4</sup> over a random sample chosen from the local data set  $\mathcal{D}_r$ . For  $f: \mathbb{R}^d \to \mathbb{R}$ , we denote  $\mathbf{x}^* := \arg\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$  and  $f^* := f(\mathbf{x}^*)$ . The distributed nodes perform computations and provide updates to the master node that is responsible for aggregation and model update. We develop *Qsparse-local-SGD*, a distributed SGD composing gradient quantization and explicit sparsification (e.g.,  $\text{Top}_k$  components), along with local iterations. We develop the algorithms and analysis for both synchronous as well as asynchronous operations, in which workers can communicate with the master at arbitrary time intervals. To the best of our knowledge, these are the first algorithms which combine quantization, aggressive sparsification, and local computations for distributed optimization.

<sup>&</sup>lt;sup>2</sup>In our asynchronous model, the distributed nodes' iterates evolve at the same rate, but update the gradients at arbitrary times; see Section 4 for more details.

<sup>&</sup>lt;sup>3</sup>At the completion of our work, we recently found that in parallel to our work [15] examined use of sign-SGD quantization, *without sparsification* for the centralized model. Another recent work in [16] studies the decentralized case with sparsification for strongly convex function. Our work, developed independent of these works, uses quantization, sparsification and local computations for the distributed case with local computations for both non-convex and strongly convex objectives.

<sup>&</sup>lt;sup>4</sup>Our setup can also handle different local functional forms, beyond dependence on the local data set  $\mathcal{D}_r$ , which is not explicitly written for notational simplicity.

Our main theoretical results are the convergence analysis of *Qsparse-local-SGD* for both (smooth) non-convex objectives as well as for the strongly convex case. See Theorem 1, 2 for the synchronous case, as well as Theorem 3, 4, for the asynchronous operation. Our analysis also demonstrates natural gains in convergence that distributed, mini-batch operation affords, and has convergence similar to vanilla SGD with local iterations (see Corollary 1, 2), for both the non-convex case (with convergence rate  $\sim 1/\sqrt{T}$  for fixed learning rate) as well as the strongly convex case (with convergence rate  $\sim 1/T$ , for diminishing learning rate), demonstrating that quantizing and sparsifying the gradient, even after local iterations asymptotically yields an almost "free" communication efficiency gain (also observed numerically in Section 5 non-asymptotically). The numerical results on ImageNet dataset implemented for a ResNet-50 architecture demonstrates that one can get significant communication savings, while retaining equivalent state-of-the art performance with a small penalty in final accuracy.

Unlike previous works, Qsparse-local-SGD stores the compression error of the net  $local\ update$ , which is a sum of at most H gradient steps and the historical error, in the local memory. From literature [4,30], we know that methods with error compensation work only when the evolution of the error is controlled. The combination of quantization, sparsification, and local computations poses several challenges for theoretical analysis, including (i) the analysis of impact of local iterations on the evolution of the error due to quantization and sparsification, as well as the deviation of local iterates (see Lemma 3, 4, 8, 9) (ii) asynchronous updates together with distribution compression using operators which satisfy Definition 3, including our composed (Qsparse) operators. (see Lemma 11-14 in appendix). Another useful technical observation is that the composition of a quantizer and a sparsifier results in a compression operator (Lemma 1, 2); see Appendix A for proofs on the same.

We provide additional results in the appendices as part of the supplementary material. These include results on the asymptotic analysis for non-convex objectives in Theorem 5, 8 along with precise statements of the convergence guarantees for the asynchronous operation Theorem 6, 7 and numerics for the convex case for multi-class logistic classification on *MNIST* [19] dataset in Appendix D, for both synchronous and asynchronous operations.

We believe that our approach for combining different forms of compression and local computations can be extended to the decentralized case, where nodes are connected over an arbitrary graph, building on the ideas from [15, 35]. Our numerics also incorporate momentum acceleration, whose analysis is a topic for future research, for example incorporating ideas from [42].

**Organization.** In Section 2, we demonstrate that composing certain classes of quantization with sparsification satisfies a certain regularity condition that is needed for several convergence proofs for our algorithms. We describe the synchronous implementation of *Qsparse-local-SGD* in Section 3, and outline the main convergence results for it in Section 3.1, briefly giving the proof ideas in Section 3.2. We describe our asynchronous implementation of *Qsparse-local-SGD* and provide the theoretical convergence results in Section 4. The experimental results are given in Section 5. Many of the proof details and additional results are given in the appendices provided with the supplementary material.

## 2 Composition of Quantization and Sparsification

In this section, we consider composition of two different techniques used in the literature for mitigating the communication bottleneck in distributed optimization, namely, quantization and sparsification. In quantization, we reduce precision of the gradient vector by mapping each of its components by a deterministic [7,15] or randomized [3,33,39,44] map to a finite number of quantization levels. In sparsification, we sparsify the gradients vector before using it to update the parameter vector, by taking its  $\mathrm{Top}_k$  components or choosing k components uniformly at random, denoted by  $\mathrm{Rand}_k$ , [30]. **Definition 1** (Randomized Quantizer [3,33,39,44]). We say that  $Q_s:\mathbb{R}^d\to\mathbb{R}^d$  is a randomized quantizer with s quantization levels, if the following holds for every  $\mathbf{x}\in\mathbb{R}^d$ : (i)  $\mathbb{E}_Q[Q_s(\mathbf{x})]=\mathbf{x}$ ; (ii)  $\mathbb{E}_Q[\|Q_s(\mathbf{x})\|^2] \leq (1+\beta_{d,s})\|\mathbf{x}\|^2$ , where  $\beta_{d,s}>0$  could be a function of d and s. Here expectation is taken over the randomness of  $Q_s$ .

Examples of randomized quantizers include (i) QSGD [3,39], which independently quantizes components of  $\mathbf{x} \in \mathbb{R}^d$  into s levels, with  $\beta_{d,s} = \min(\frac{d}{s^2}, \frac{\sqrt{d}}{s})$ ; (ii) Stochastic s-level Quantization [33,44], which independently quantizes every component of  $\mathbf{x} \in \mathbb{R}^d$  into s levels between  $\operatorname{argmax}_i x_i$  and  $\operatorname{argmin}_i x_i$ , with  $\beta_{d,s} = \frac{d}{2s^2}$ ; and (iii) Stochastic Rotated Quantization [33], which is a stochastic quantization, preprocessed by a random rotation, with  $\beta_{d,s} = \frac{2\log_2(2d)}{s^2}$ .

Instead of quantizing randomly into s levels, we can take a deterministic approach and round off to the nearest level. In particular, we can just take the sign, which has shown promise in [7, 27, 32].

**Definition 2** (Deterministic Sign Quantizer [7, 15]). A deterministic quantizer  $Sign : \mathbb{R}^d \to \{+1, -1\}^d$  is defined as follows: for every vector  $\mathbf{x} \in \mathbb{R}^d$ ,  $i \in [d]$ , the i'th component of  $Sign(\mathbf{x})$  is defined as  $\mathbb{1}\{x_i \geq 0\} - \mathbb{1}\{x_i < 0\}$ .

As mentioned above, we consider two important examples of sparsification operators:  $\operatorname{Top}_k$  and  $\operatorname{Rand}_k$ , For any  $\mathbf{x} \in \mathbb{R}^d$ ,  $\operatorname{Top}_k(\mathbf{x})$  is equal to a d-length vector, which has at most k non-zero components whose indices correspond to the indices of the largest k components (in absolute value) of  $\mathbf{x}$ . Similarly,  $\operatorname{Rand}_k(\mathbf{x})$  is a d-length (random) vector, which is obtained by selecting k components of  $\mathbf{x}$  uniformly at random. Both of these satisfy a so-called "compression" property as defined below, with  $\gamma = k/d$  [30]. Few other examples of such operators can be found in [30].

**Definition 3** (Sparsification [30]). A (randomized) function  $Comp_k : \mathbb{R}^d \to \mathbb{R}^d$  is called a compression operator, if there exists a constant  $\gamma \in (0,1]$  (that may depend on k and d), such that for every  $\mathbf{x} \in \mathbb{R}^d$ , we have  $\mathbb{E}_C[\|\mathbf{x} - Comp_k(\mathbf{x})\|_2^2] \leq (1-\gamma)\|\mathbf{x}\|_2^2$ , where expectation is taken over  $Comp_k$ .

We can apply different compression operators to different coordinates of a vector, and the resulting operator is also a compression operator; see Corollary 3 in Appendix A. As an application, in the case of training neural networks, we can apply different compression operators to different layers.

**Composition of Quantization and Sparsification.** Now we show that we can compose deterministic/randomized quantizers with sparsifiers and the resulting operator is a compression operator. Proofs are given in Appendix A.

**Lemma 1** (Composing sparsification with stochastic quantization). Let  $Comp_k \in \{ Top_k, Rand_k \}$ . Let  $Q_s : \mathbb{R}^d \to \mathbb{R}^d$  be a stochastic quantizer with parameter s that satisfies Definition 1. Let  $Q_sComp_k : \mathbb{R}^d \to \mathbb{R}^d$  be defined as  $Q_sComp_k(\mathbf{x}) := Q_s(Comp_k(\mathbf{x}))$  for every  $\mathbf{x} \in \mathbb{R}^d$ . Then  $\frac{Q_sComp_k(\mathbf{x})}{1+\beta_{k,s}}$  is a compression operator with the compression coefficient being equal to  $\gamma = \frac{k}{d(1+\beta_{k,s})}$ .

**Lemma 2** (Composing sparsification with deterministic quantization). Let  $Comp_k \in \{ Top_k, Rand_k \}$ . Let  $SignComp_k : \mathbb{R}^d \to \mathbb{R}^d$  be defined as follows: for every  $\mathbf{x} \in \mathbb{R}^d$ , the i'th component of  $SignComp_k(\mathbf{x})$  is equal to  $\mathbb{1}\{x_i \geq 0\} - \mathbb{1}\{x_i < 0\}$ , if the i'th component is chosen in defining  $Comp_k$ , otherwise, it is equal to 0. Then  $\frac{\|Comp_k(\mathbf{x})\|_1 \ SignComp_k(\mathbf{x})}{k}$  is a compression operator<sup>5</sup> with the compression coefficient being equal to  $\gamma = \max \left\{ \frac{1}{d}, \frac{k}{d} \left( \frac{\|Comp_k(\mathbf{x})\|_1}{\sqrt{d}\|Comp_k(\mathbf{x})\|_2} \right)^2 \right\}$ .

## 3 Qsparse-local-SGD

Let  $\mathcal{I}_T^{(r)}\subseteq [T]:=\{1,\ldots,T\}$  with  $T\in\mathcal{I}_T^{(r)}$  denote a set of indices for which worker  $r\in [R]$  synchronizes with the master. In a synchronous setting,  $\mathcal{I}_T^{(r)}$  is same for all the workers. Let  $\mathcal{I}_T:=\mathcal{I}_T^{(r)}$  for any  $r\in [R]$ . Every worker  $r\in [R]$  maintains a local parameter  $\widehat{\mathbf{x}}_t^{(r)}$  which is updated in each iteration t, using the stochastic gradient  $\nabla f_{i_t^{(r)}}\left(\widehat{\mathbf{x}}_t^{(r)}\right)$ , where  $i_t^{(r)}$  is a mini-batch of size t0 sampled uniformly in t2. If  $t\in \mathcal{I}_T$ , the sparsified error-compensated update t3 computed on the net progress made since the last synchronization is sent to the master node, and updates its local memory t4. Upon receiving t7 is from every worker, master aggregates them, updates the global parameter vector, and sends the new model t4 to all the workers; upon receiving which, they set their local parameter vector t6. Our algorithm is summarized in Algorithm 1.

### 3.1 Main Results for Synchronous Operation

All results in this paper use the following two standard assumptions. (i) **Smoothness:** The local function  $f^{(r)}: \mathbb{R}^d \to \mathbb{R}$  at each worker  $r \in [R]$  is L-smooth, i.e., for every  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ , we have  $f^{(r)}(\mathbf{y}) \leq f^{(r)}(\mathbf{x}) + \langle \nabla f^{(r)}(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} ||\mathbf{y} - \mathbf{x}||^2$ . (ii) **Bounded second moment:** For every

The analysis for general p-norm, i.e.  $\frac{\|Comp_k(\mathbf{x})\|_p \ SignComp_k(\mathbf{x})}{k}$ , for any  $p \in \mathbb{Z}_+$  is provided in Appendix A.

## Algorithm 1 Qsparse-local-SGD

```
1: Initialize \mathbf{x}_0 = \widehat{\mathbf{x}}_0^{(r)} = m_0^{(r)}, \ \forall r \in [R]. Suppose \eta_t follows a certain learning rate schedule. 2: for t = 0 to T - 1 do
               On Workers:
              \begin{aligned} & \textbf{for } r = 1 \textbf{ to } R \textbf{ do} \\ & \widehat{\mathbf{x}}_{t+\frac{1}{2}}^{(r)} \leftarrow \widehat{\mathbf{x}}_t^{(r)} - \eta_t \nabla f_{i_t^{(r)}} \left( \widehat{\mathbf{x}}_t^{(r)} \right); i_t^{(r)} \text{ is a mini-batch of size } b \text{ sampled uniformly in } \mathcal{D}_r \\ & \textbf{ if } t + 1 \notin \mathcal{I}_T \textbf{ then} \end{aligned}
  5:
  6:
                          \mathbf{x}_{t+1} \leftarrow \mathbf{x}_t, m_{t+1}^{(r)} \leftarrow m_t^{(r)} \text{ and } \widehat{\mathbf{x}}_{t+1}^{(r)} \leftarrow \widehat{\mathbf{x}}_{t+\frac{1}{2}}^{(r)}
  7:
                    else g_t^{(r)} \leftarrow Q \, Comp_k \left( m_t^{(r)} + \mathbf{x}_t - \widehat{\mathbf{x}}_{t+\frac{1}{2}}^{(r)} \right), \text{ send } g_t^{(r)} \text{ to the master.} m_{t+1}^{(r)} \leftarrow m_t^{(r)} + \mathbf{x}_t - \widehat{\mathbf{x}}_{t+\frac{1}{2}}^{(r)} - g_t^{(r)}
  8:
 9:
10:
                           Receive \mathbf{x}_{t+1} from the master and set \widehat{\mathbf{x}}_{t+1}^{(r)} \leftarrow \mathbf{x}_{t+1}
11:
12:
                     end if
13:
                end for
                At Master:
14:
15:
               if t+1 \notin \mathcal{I}_T then
16:
17:
                     Receive g_t^{(r)} from R workers and compute \mathbf{x}_{t+1} = \mathbf{x}_t - \frac{1}{R} \sum_{r=1}^{R} g_t^{(r)}
18:
19:
20:
21: end for
22: Comment: Note that \widehat{\mathbf{x}}_{t+\frac{1}{2}}^{(r)} is used to denote an intermediate variable between iterations t and t+1.
```

 $\widehat{\mathbf{x}}_t^{(r)} \in \mathbb{R}^d, r \in [R], t \in [T]$ , we have  $\underset{i \sim \mathcal{D}_r}{\mathbb{E}} [\|\nabla f_i(\widehat{\mathbf{x}}_t^{(r)})\|^2] \leq G^2$ , for some constant  $G < \infty$ . This is a standard assumption in [4, 12, 16, 23, 25, 26, 29–31, 43]. Relaxation of the uniform boundedness of the gradient allowing arbitrarily different gradients of local functions in heterogenous settings as done for SGD in [24, 37] is left as future work. This also imposes a **bound on the variance**:  $\underset{i \sim \mathcal{D}_r}{\mathbb{E}} [\|\nabla f_i(\widehat{\mathbf{x}}_t^{(r)}) - \nabla f^{(r)}(\widehat{\mathbf{x}}_t^{(r)})\|^2] \leq \sigma_r^2, \text{ where } \sigma_r^2 \leq G^2 \text{ for every } r \in [R]. \text{ To state our results, we need the following definition from [31].}$ 

**Definition 4** (Gap [31]). Let  $\mathcal{I}_T = \{t_0, t_1, \dots, t_k\}$ , where  $t_i < t_{i+1}$  for  $i = 0, 1, \dots, k-1$ . The gap of  $\mathcal{I}_T$  is defined as  $gap(\mathcal{I}_T) := \max_{i \in [k]} \{(t_i - t_{i-1})\}$ , which is equal to the maximum difference between any two consecutive synchronization indices.

We leverage the perturbed iterate analysis as in [21, 30] to provide convergence guarantees for *Qsparse-local-SGD*. Under assumptions (i) and (ii), the following theorems hold when Algorithm 1 is run with any compression operator (including our composed operators).

**Theorem 1** (Convergence in the smooth (non-convex) case with fixed learning rate). Let  $f^{(r)}(\mathbf{x})$  be L-smooth for every  $i \in [R]$ . Let  $QComp_k : \mathbb{R}^d \to \mathbb{R}^d$  be a compression operator whose compression coefficient is equal to  $\gamma \in (0,1]$ . Let  $\{\widehat{\mathbf{x}}_t^{(r)}\}_{t=0}^{T-1}$  be generated according to Algorithm 1 with  $QComp_k$ , for step sizes  $\eta = \frac{\widehat{C}}{\sqrt{T}}$  (where  $\widehat{C}$  is a constant such that  $\frac{\widehat{C}}{\sqrt{T}} \leq \frac{1}{2L}$ ) and  $gap(\mathcal{I}_T) \leq H$ . Then we have

$$\mathbb{E}\|\nabla f(\mathbf{z}_T)\|^2 \le \left(\frac{\mathbb{E}[f(\mathbf{x}_0)] - f^*}{\widehat{C}} + \widehat{C}L\left(\frac{\sum_{r=1}^R \sigma_r^2}{bR^2}\right)\right) \frac{4}{\sqrt{T}} + 8\left(4\frac{(1-\gamma^2)}{\gamma^2} + 1\right) \frac{\widehat{C}^2L^2G^2H^2}{T}.\tag{1}$$

Here  $\mathbf{z}_T$  is a random variable which samples a previous parameter  $\widehat{\mathbf{x}}_t^{(r)}$  with probability 1/RT.

**Corollary 1.** Let  $\mathbb{E}[f(\mathbf{x}_0)] - f^* \leq J^2$ , where  $J < \infty$  is a constant,  $\sigma_{max} = \max_{r \in [R]} \sigma_r$ , and  $\widehat{C}^2 = \frac{bR(\mathbb{E}[f(\mathbf{x}_0)] - f^*)}{\sigma_{max}^2 L}$ , we have

$$\mathbb{E}\|\nabla f(\mathbf{z}_T)\|^2 \le \mathcal{O}\left(\frac{J\sigma_{max}}{\sqrt{bRT}}\right) + \mathcal{O}\left(\frac{J^2bRG^2H^2}{\sigma_{max}^2\gamma^2T}\right). \tag{2}$$

<sup>&</sup>lt;sup>6</sup>Even classical SGD requires knowing an upper bound on  $\|\mathbf{x}_0 - \mathbf{x}^*\|$  in order to choose the learning rate. Smoothness of f translates this to the difference of the function values.

In order to ensure that the compression does not affect the dominating terms while converging at a rate of  $\mathcal{O}\left(1/\sqrt{bRT}\right)$ , we would require  $H = \mathcal{O}\left(\gamma T^{1/4}/(bR)^{3/4}\right)$ .

Theorem 1 is proved in Appendix B and provides non-asymptotic guarantees, where we observe that compression does not affect the first order term. The corresponding asymptotic result (with decaying learning rate), with a convergence rate of  $\mathcal{O}(\frac{1}{\log T})$ , is provided in Theorem 5 in Appendix B.

**Theorem 2** (Convergence in the smooth and strongly convex case with a decaying learning rate). Let  $f^{(r)}(\mathbf{x})$  be L-smooth and  $\mu$ -strongly convex. Let  $QComp_k: \mathbb{R}^d \to \mathbb{R}^d$  be a compression operator whose compression coefficient is equal to  $\gamma \in (0,1]$ . Let  $\{\widehat{\mathbf{x}}_t^{(r)}\}_{t=0}^{T-1}$  be generated according to Algorithm 1 with  $QComp_k$ , for step sizes  $\eta_t = \frac{8}{\mu(a+t)}$  with  $gap(\mathcal{I}_T) \leq H$ , where a > 1 is such that we have  $a \geq \max\{\frac{4H}{\gamma}, 32\kappa, H\}$ ,  $\kappa = \frac{L}{\mu}$ . Then the following holds

$$\mathbb{E}[f(\overline{\mathbf{x}}_T)] - f^* \le \frac{La^3}{4S_T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 + \frac{8LT(T + 2a)}{\mu^2 S_T} A + \frac{128LT}{\mu^3 S_T} B.$$
 (3)

Here (i) 
$$A = \frac{\sum_{r=1}^{R} \sigma_r^2}{bR^2}$$
,  $B = 4\left(\left(\frac{3\mu}{2} + 3L\right)\frac{CG^2H^2}{\gamma^2} + 3L^2G^2H^2\right)$ , where  $C \geq \frac{4a\gamma(1-\gamma^2)}{a\gamma-4H}$ ; (ii)  $\overline{\mathbf{x}}_T := \frac{1}{S_T}\sum_{t=0}^{T-1} \left[w_t\left(\frac{1}{R}\sum_{r=1}^{R}\widehat{\mathbf{x}}_t^{(r)}\right)\right]$ , where  $w_t = (a+t)^2$ ; and (iii)  $S_T = \sum_{t=0}^{T-1} w_t \geq \frac{T^3}{3}$ .

**Corollary 2.** For  $a > \max\{\frac{4H}{\gamma}, 32\kappa, H\}$ ,  $\sigma_{max} = \max_{r \in [R]} \sigma_r$ , and using  $\mathbb{E} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \leq \frac{4G^2}{\mu^2}$  from Lemma 2 in [25], we have

$$\mathbb{E}[f(\overline{\mathbf{x}}_T)] - f^* \le \mathcal{O}\left(\frac{G^2 H^3}{\mu^2 \gamma^3 T^3}\right) + \mathcal{O}\left(\frac{\sigma_{max}^2}{\mu^2 b R T} + \frac{H \sigma_{max}^2}{\mu^2 b R \gamma T^2}\right) + \mathcal{O}\left(\frac{G^2 H^2}{\mu^3 \gamma^2 T^2}\right). \tag{4}$$

In order to ensure that the compression does not affect the dominating terms while converging at a rate of  $\mathcal{O}(1/(bRT))$ , we would require  $H = \mathcal{O}\left(\gamma\sqrt{T/(bR)}\right)$ .

Theorem 2 has been proved in Appendix B. For no compression and only local computations, i.e., for  $\gamma=1$ , and under the same assumptions, we recover/generalize a few recent results from literature with similar convergence rates: (i) We recover [43, Theorem 1], which is for non-convex case; (ii) We generalize [31, Theorem 2.2], which is for a strongly convex case and requires that each worker has identical datasets, to the distributed case. We emphasize that unlike [31,43], which only consider local computation, we combine quantization and sparsification with local computation, which poses several technical challenges (e.g., see proofs of Lemma 3, 4,7 in Appendix B).

#### 3.2 Proof Outlines

Maintain virtual sequences for every worker

$$\widetilde{\mathbf{x}}_{0}^{(r)} := \widehat{\mathbf{x}}_{0}^{(r)} \quad \text{and} \quad \widetilde{\mathbf{x}}_{t+1}^{(r)} := \widetilde{\mathbf{x}}_{t}^{(r)} - \eta_{t} \nabla f_{i}^{(r)} \left(\widehat{\mathbf{x}}_{t}^{(r)}\right)$$
 (5)

Define (i) 
$$\mathbf{p}_t := \frac{1}{R} \sum_{r=1}^R \nabla f_{i_t^{(r)}} \left( \widehat{\mathbf{x}}_t^{(r)} \right), \quad \overline{\mathbf{p}}_t := \mathbb{E}_{i_t}[\mathbf{p}_t] = \frac{1}{R} \sum_{r=1}^R \nabla f^{(r)} \left( \widehat{\mathbf{x}}_t^{(r)} \right);$$

and (ii) 
$$\widetilde{\mathbf{x}}_{t+1} := \frac{1}{R} \sum_{r=1}^{R} \widetilde{\mathbf{x}}_{t+1}^{(r)} = \widetilde{\mathbf{x}}_{t} - \eta_{t} \mathbf{p}_{t}, \quad \widehat{\mathbf{x}}_{t} := \frac{1}{R} \sum_{r=1}^{R} \widehat{\mathbf{x}}_{t}^{(r)}.$$

Proof outline of Theorem 1. Since f is L-smooth, we have  $f(\widetilde{\mathbf{x}}_{t+1}) - f(\widetilde{\mathbf{x}}_t) \leq -\eta_t \langle \nabla f(\widetilde{\mathbf{x}}_t), \mathbf{p}_t \rangle + \frac{\eta_t^2 L}{2} \|\mathbf{p}_t\|^2$ . With some algebraic manipulations provided in Appendix B, for  $\eta_t \leq 1/2L$ , we arrive at

$$\frac{\eta_t}{4R} \sum_{r=1}^R \mathbb{E} \|\nabla f(\widehat{\mathbf{x}}_t^{(r)})\|^2 \le \mathbb{E}[f(\widetilde{\mathbf{x}}_t)] - \mathbb{E}[f(\widetilde{\mathbf{x}}_{t+1})] + \eta_t^2 L \mathbb{E} \|\mathbf{p}_t - \overline{\mathbf{p}}_t\|^2 + 2\eta_t L^2 \mathbb{E} \|\widetilde{\mathbf{x}}_t - \widehat{\mathbf{x}}_t\|^2 + 2\eta_t L^2 \frac{1}{R} \sum_{r=1}^R \mathbb{E} \|\widehat{\mathbf{x}}_t - \widehat{\mathbf{x}}_t^{(r)}\|^2.$$
(6)

Under Assumptions 1 and 2, we have  $\mathbb{E}\|\mathbf{p}_t - \overline{\mathbf{p}}_t\|^2 \leq \frac{\sum_{r=1}^R \sigma_r^2}{bR^2}$ . To bound  $\mathbb{E}\|\widetilde{\mathbf{x}}_t - \widehat{\mathbf{x}}_t\|^2$  in (6), we first show (in Lemma 7 in Appendix B) that  $\widehat{\mathbf{x}}_t - \widetilde{\mathbf{x}}_t = \frac{1}{R}\sum_{r=1}^R m_t^{(r)}$ , i.e., the difference of the true and the virtual parameter vectors is equal to the average memory, and then we bound the local memory at each worker  $r \in [R]$  below.

<sup>&</sup>lt;sup>7</sup>Here we characterize the reduction in communication that can be afforded, however for a constant H we get the same rate of convergence after  $T = \Omega\left((bR)^3/\gamma^4\right)$ . Analogous statements hold for Theorem 2-4.

**Lemma 3** (Bounded Memory). For  $\eta_t = \eta$ ,  $gap(\mathcal{I}_T) \leq H$ , we have for every  $t \in \mathbb{Z}^+$  that

$$\mathbb{E}\|m_t^{(r)}\|^2 \le 4\frac{\eta^2(1-\gamma^2)}{\gamma^2}H^2G^2. \tag{7}$$

Using Lemma 3, we get  $\mathbb{E}\|\widehat{\mathbf{x}}_t-\widehat{\mathbf{x}}_t\|^2 \leq \frac{1}{R}\sum_{r=1}^R \mathbb{E}\|m_t^{(r)}\|^2 \leq 4\frac{\eta^2(1-\gamma^2)}{\gamma^2}H^2G^2$ . We can bound the last term of (6) as  $\frac{1}{R}\sum_{r=1}^R \mathbb{E}\|\widehat{\mathbf{x}}_t-\widehat{\mathbf{x}}_t^{(r)}\|^2 \leq \eta^2G^2H^2$  in Lemma 9 in Appendix B. Putting them back in (6), performing a telescopic sum from t=0 to T-1, and then taking an average over time, we get

$$\frac{1}{RT} \sum_{t=0}^{T-1} \sum_{r=1}^{R} \mathbb{E} \|\nabla f(\widehat{\mathbf{x}}_t^{(r)})\|^2 \leq \frac{4 \left(\mathbb{E}[f(\widehat{\mathbf{x}}_0)] - f^*\right)}{\eta T} + \frac{4\eta L}{bR^2} \sum_{r=1}^{R} \sigma_r^2 + 32 \frac{\eta^2 (1-\gamma^2)}{\gamma^2} L^2 G^2 H^2 + 8\eta^2 L^2 G^2 H^2.$$

By letting  $\eta = \hat{C}/\sqrt{T}$ , where  $\hat{C}$  is a constant such that  $\frac{\hat{C}}{\sqrt{T}} \leq \frac{1}{2L}$ , we arrive at Theorem 1.

Proof outline of Theorem 2. Using the definition of virtual sequences (5), we have  $\|\widetilde{\mathbf{x}}_{t+1} - \mathbf{x}^*\|^2 = \|\widetilde{\mathbf{x}}_t - \mathbf{x}^* - \eta_t \overline{\mathbf{p}}_t\|^2 + \eta_t^2 \|\mathbf{p}_t - \overline{\mathbf{p}}_t\|^2 - 2\eta_t \langle \widetilde{\mathbf{x}}_t - \mathbf{x}^* - \eta_t \overline{\mathbf{p}}_t, \mathbf{p}_t - \overline{\mathbf{p}}_t \rangle$ . With some algebraic manipulations provided in Appendix B, for  $\eta_t \leq 1/4L$  and letting  $e_t = \mathbb{E}[f(\widehat{\mathbf{x}}_t)] - f^*$ , we get

$$\mathbb{E}\|\widetilde{\mathbf{x}}_{t+1} - \mathbf{x}^*\|^2 \le \left(1 - \frac{\mu \eta_t}{2}\right) \mathbb{E}\|\widetilde{\mathbf{x}}_t - \mathbf{x}^*\|^2 - \frac{\eta_t \mu}{2L} e_t + \eta_t \left(\frac{3\mu}{2} + 3L\right) \mathbb{E}\|\widehat{\mathbf{x}}_t - \widetilde{\mathbf{x}}_t\|^2 + \frac{3\eta_t L}{R} \sum_{r=1}^R \mathbb{E}\|\widehat{\mathbf{x}}_t - \widehat{\mathbf{x}}_t^{(r)}\|^2 + \eta_t^2 \frac{\sum_{r=1}^R \sigma_r^2}{bR^2}.$$
 (8)

To bound the 3rd term on the RHS of (63), first we note that  $\hat{\mathbf{x}}_t - \tilde{\mathbf{x}}_t = \frac{1}{R} \sum_{r=1}^R m_t^{(r)}$ , and then we bound the local memory at each worker  $r \in [R]$  below.

**Lemma 4** (Memory Contraction). For  $a > {}^{4H}/\gamma$ ,  $\eta_t = \xi/a + t$ ,  $gap(\mathcal{I}_T) \leq H$ , there exists a  $C \geq \frac{4a\gamma(1-\gamma^2)}{a\gamma-4H}$  such that the following holds for every  $t \in \mathbb{Z}^+$ 

$$\mathbb{E}\|m_t^{(r)}\|^2 \le 4\frac{\eta_t^2}{2}CH^2G^2. \tag{9}$$

A proof of Lemma 4 is provided in Appendix B and is technically more involved than the proof of Lemma 3. This complication arises because of the decaying learning rate, combined with compression and local computation. We can bound the penultimate term on the RHS of (63) as  $\frac{1}{R}\sum_{r=1}^{R}\mathbb{E}\|\widehat{\mathbf{x}}_t-\widehat{\mathbf{x}}_t^{(r)}\|^2 \leq 4\eta_t^2G^2H^2.$  This can be shown along the lines of the proof of [31, Lemma 3.3] and we show it in Lemma 8 in Appendix B. Substituting all these in (63) gives

$$\mathbb{E}\|\widetilde{\mathbf{x}}_{t+1} - \mathbf{x}^*\|^2 \le \left(1 - \frac{\mu \eta_t}{2}\right) \mathbb{E}\|\widetilde{\mathbf{x}}_t - \mathbf{x}^*\|^2 - \frac{\mu \eta_t}{2L} e_t + \eta_t \left(\frac{3\mu}{2} + 3L\right) C \frac{4\eta_t^2}{\gamma^2} G^2 H^2 + (3\eta_t L) 4\eta_t^2 L G^2 H^2 + \eta_t^2 \frac{\sum_{r=1}^R \sigma_r^2}{2r}.$$
(10)

Since (10) is a contracting recurrence relation, with some calculation done in Appendix B, we complete the proof of Theorem 2.

## 4 Asynchronous Osparse-local-SGD

We propose and analyze a particular form of asynchronous operation where the workers synchronize with the master at arbitrary times decided locally or by master picking a subset of nodes as in federated learning [17,22]. However, the local iterates evolve at the same rate, i.e. each worker takes the same number of steps per unit time according to a global clock. The asynchrony is therefore that updates occur after different number of local iterations but the local iterations are synchronous with respect to the global clock.<sup>8</sup>

In this asynchronous setting,  $\mathcal{I}_T^{(r)}$ 's may be different for different workers. However, we assume that  $gap(\mathcal{I}_T^{(r)}) \leq H$  holds for every  $r \in [R]$ , which means that there is a uniform bound on the maximum delay in each worker's update times. The algorithmic difference from Algorithm 1 is that, in this case, a subset of workers (including a single worker) can send their updates to the master at their synchronization time steps; master aggregates them, updates the global parameter vector, and sends that only to those workers. Our algorithm is summarized in Algorithm 2 in Appendix C. We give the simplified expressions of our main results below; more precise results are in Appendix C.

<sup>&</sup>lt;sup>8</sup>This is different from asynchronous algorithms studied for stragglers [26,41], where only one gradient step is taken but occurs at different times due to delays.

Theorem 3 (Convergence in the smooth non-convex case with fixed learning rate). Under the same conditions as in Theorem 1 with  $gap(\mathcal{I}_T^{(r)}) \leq H$ , if  $\{\widehat{\mathbf{x}}_t^{(r)}\}_{t=0}^{T-1}$  is generated according to Algorithm 2, the following holds, where  $\mathbb{E}[f(\mathbf{x}_0)] - f^* \leq J^2$ ,  $\sigma_{max} = \max_{r \in [R]} \sigma_r$ , and  $\widehat{C}^2 = \sigma_{max}$  $bR(\mathbb{E}[f(\mathbf{x}_0)] - f^*) / \sigma_{max}^2.$   $\mathbb{E}\|\nabla f(\mathbf{z}_T)\|^2 \le \mathcal{O}\left(\frac{J\sigma_{max}}{\sqrt{bRT}}\right) + \mathcal{O}\left(\frac{J^2bRG^2}{\sigma_{max}^2\gamma^2T}(H^2 + H^4)\right).$ 

(11)

where  $\mathbf{z}_T$  is a random variable which samples a previous parameter  $\widehat{\mathbf{x}}_t^{(r)}$  with probability 1/RT. In order to ensure that the compression does not affect the dominating terms while converging at a rate of  $\mathcal{O}\left(1/\sqrt{bRT}\right)$ , we would require  $H = \mathcal{O}\left(\sqrt{\gamma}T^{1/8}/(bR)^{3/8}\right)$ .

We give a precise result in Theorem 6 in Appendix C. Note that Theorem 3 provides non-asymptotic guarantees, where compression is almost for "free". The corresponding asymptotic result with decaying learning rate, with a convergence rate of  $\mathcal{O}(\frac{1}{\log T})$ , is provided in Theorem 8 in Appendix C.

Theorem 4 (Convergence in the smooth and strongly convex case with decaying learning rate). Under the same conditions as in Theorem 2 with  $gap(\mathcal{I}_T^{(r)}) \leq H$ ,  $a > \max\{4H/\gamma, 32\kappa, H\}$ ,  $\sigma_{max} = 0$  $\max_{r \in [R]} \sigma_r$ , if  $\{\widehat{\mathbf{x}}_t^{(r)}\}_{t=0}^{T-1}$  is generated according to Algorithm 2, the following holds:

$$\mathbb{E}[f(\overline{\mathbf{x}}_T)] - f^* \le \mathcal{O}\left(\frac{G^2 H^3}{\mu^2 \gamma^3 T^3}\right) + \mathcal{O}\left(\frac{\sigma_{max}^2}{\mu^2 b R T} + \frac{H \sigma_{max}^2}{\mu^2 b R \gamma T^2}\right) + \mathcal{O}\left(\frac{G^2}{\mu^3 \gamma^2 T^2} (H^2 + H^4)\right). \tag{12}$$

where  $\overline{\mathbf{x}}_T$ ,  $S_T$  are as defined in Theorem 2. To ensure that the compression does not affect the dominating terms while converging at a rate of  $\mathcal{O}(1/(bRT))$ , we would require  $H = \mathcal{O}(\sqrt{\gamma}(T/(bR))^{1/4})$ .

We give a more precise result in Theorem 7 in Appendix C. If  $\mathcal{I}_T^{(r)}$ 's are the same for all the workers, then one would ideally require that the bounds on H in the asynchronous setting reduce to the bounds on H in the synchronous setting. This is not happening, as our bounds in the asynchronous setting are for the worst case scenario – they hold as long as  $gap(\mathcal{I}_T^{(r)}) \leq H$ , for every  $r \in [R]$ .

## 4.1 Proof Outlines

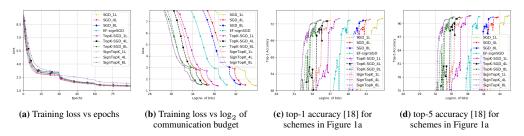
Our proofs of these results follow the same outlines of the corresponding proofs in the synchronous setting, but some technical details change significantly. This is because, in our asynchronous setting, workers are allowed to update the global parameter vector in between two consecutive synchronization time steps of other workers. For example, unlike the synchronous setting,  $\widehat{\mathbf{x}}_t - \widetilde{\mathbf{x}}_t = \frac{1}{R} \sum_{r=1}^R m_t^{(r)}$  does not hold here; however, we can show that  $\widehat{\mathbf{x}}_t - \widetilde{\mathbf{x}}_t$  is equal to the sum of  $\frac{1}{R} \sum_{r=1}^R m_t^{(r)}$  and an additional term, which leads to potentially a weaker bound  $\mathbb{E} \|\widehat{\mathbf{x}}_t - \widetilde{\mathbf{x}}_t\|^2 \leq \mathcal{O}\left(n_t^2/\gamma^2 G^2(H^2 + H^4)\right)$ (vs.  $\mathcal{O}\left(\eta_t^2/\gamma^2G^2H^2\right)$  for the synchronous setting), proved in Lemma 13-14 in Appendix C. Similarly, the proof of the average true sequence being close to the virtual sequence requires carefully chosen reference points on the global parameter sequence lying within bounded steps of the local parameters. We show a bound on  $\frac{1}{R}\sum_{r=1}^{R}\mathbb{E}\|\widehat{\mathbf{x}}_t-\widehat{\mathbf{x}}_t^{(r)}\|^2 \leq \mathcal{O}(\eta_t^2G^2(H^2+H^4/\gamma^2))$ , which is weaker than the corresponding bound  $\mathcal{O}(\eta_t^2G^2H^2)$  for the synchronous setting, in Lemma 11-12 in Appendix C.

#### 5 **Experiments**

**Experiment setup:** We train ResNet-50 [13] (which has d = 25,610,216 parameters) on ImageNet dataset, using 8 NVIDIA Tesla V100 GPUs. We use a learning rate schedule consisting of 5 epochs of linear warmup, followed by a piecewise decay of 0.1 at epochs 30, 60 and 80, with a batch size of 256 per GPU. For experiments, we focus on SGD with momentum of 0.9, applied on the local iterations of the workers. We build our compression scheme into the Horovod framework [28]. We use  $SignTop_k$  (as in Lemma 2) as our composed operator. In  $Top_k$ , we only update  $k_t = \min(d_t, 1000)$ elements per step for each tensor t, where  $d_t$  is the number of elements in the tensor. For ResNet-50 architecture, this amounts to updating a total of k = 99,400 elements per step. We also perform analogous experiments on the MNIST [19] handwritten digits dataset for softmax regression with a standard  $\ell_2$  regularizer, using the synchronous operation of *Qsparse-local-SGD* with 15 workers, and

Our implementation is available at https://github.com/karakusc/horovod/tree/qsparselocal.

a decaying learning rate as proposed in Theorem 2, the details of which are provided in Appendix D. Results: Figure 1 compares the performance of  $SignTop_k$ -SGD (which employs the 1 bit sign quantizer and the  $Top_k$  sparsifier) with error compensation (SignTopK) against (i)  $Top_k$  SGD with error compensation (TopK-SGD), (ii) SignSGD with error compensation (EF-SIGNSGD), and (iii) vanilla SGD (SGD). All of these are specializations of Qsparse-local-SGD. Furthermore, SignTopK\_hL uses a synchronization period of h; same applies for other schemes. From Figure 1a, we observe that quantization and sparsification, both individually and combined, with error compensation, has almost no penalty in terms of convergence rate, with respect to vanilla SGD. We observe that SignTopK demonstrates superior performance over EF-SIGNSGD, TopK-SGD, as well as vanilla SGD, both in terms of the required number of communicated bits for achieving a certain target loss as well as test accuracy. This is because in SignTopK, we send only 1 bit for the sign of each  $Top_k$  coordinate, along with its location. Observe that the incorporation of local iterations in Figure 1a has very little impact on the convergence rates, as compared to vanilla SGD with the same number of local iterations. Furthermore, this provides an added advantage over SignTopK, in terms of savings (by a factor of 6 to 8 times on average) in communication bits for achieving a certain target loss; see Figure 1b.



**Figure 1** Figure 1a-1d demonstrate performance gains of our of our scheme in comparison with local SGD [31], EF-SIGNSGD [15] and TopK-SGD [4, 30] in a non-convex setting for synchronous updates.

Figure 1c and Figure 1d show the top-1, and top-5 convergence rates, <sup>11</sup> respectively, with respect to the total number of bits of communication used. We observe that *Qsparse-local-SGD* combines the bit savings of the deterministic sign based operator and aggressive sparsifier, with infrequent communication; thereby, outperforming the cases where these techniques are individually used. In particular, the required number of bits to achieve the same loss or accuracy in the case of *Qsparse-local-SGD* is around 1/16 in comparison with TopK-SGD and over 1000× less than vanilla SGD.

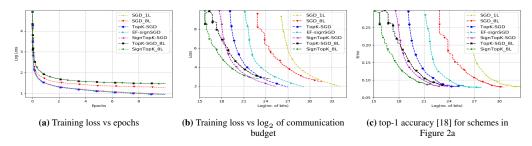


Figure 2 Figure 2a-2c demonstrate the performance gains of our scheme in a convex setting.

Figure 2b and 2c makes similar comparisons in the convex setting, and shows that for a test error approximately 0.1, Qsparse-local-SGD combines the benefits of the composed operator  $SignTop_k$ , with local computations, and needs 10-15 times less bits than TopK-SGD and  $1000 \times less$  bits than vanilla SGD. Also in Figure 2a, we observe that both TopK-SGD and SignTopK\_8L (SignTopK with 8 local iterations) converge at rates which are almost similar to that of their corresponding local SGD counterpart. Our experiments in both non-convex and convex settings verify that error compensation through memory can be used to mitigate not only the missing components from updates in previous synchronization rounds, but also explicit quantization error.

 $<sup>^{10}</sup>$ Further numerics demonstrating the performance of *Qsparse-local-SGD* for the composition of a stochastic quantizer with a sparsifier, as compared to  $SignTop_k$  and other standard baselines can be found in [6].

<sup>&</sup>lt;sup>11</sup>top-i refers to the accuracy of the top i predictions by the model from the list of possible classes; see [18].

### Acknowledgments

The authors gratefully thank Navjot Singh for his help with experiments in the early stages of this work. This work was partially supported by NSF grant #1514531, by UC-NL grant LFR-18-548554 and by Army Research Laboratory under Cooperative Agreement W911NF-17-2-0196. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Laboratory or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation here on.

## References

- [1] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, M. Kudlur, J. Levenberg, R. Monga, S. Moore, D. G. Murray, B. Steiner, P. A. Tucker, V. Vasudevan, P. Warden, M. Wicke, Y. Yu, and X. Zheng. Tensorflow: A system for large-scale machine learning. In *OSDI*, pages 265–283, 2016.
- [2] Alham Fikri Aji and Kenneth Heafield. Sparse communication for distributed gradient descent. In *EMNLP*, pages 440–445, 2017.
- [3] D. Alistarh, D. Grubic, J. Li, R. Tomioka, and M. Vojnovic. QSGD: communication-efficient SGD via gradient quantization and encoding. In *NIPS*, pages 1707–1718, 2017.
- [4] D. Alistarh, T. Hoefler, M. Johansson, N. Konstantinov, S. Khirirat, and C. Renggli. The convergence of sparsified gradient methods. In *NeurIPS*, pages 5977–5987, 2018.
- [5] Francis R. Bach and Eric Moulines. Non-asymptotic analysis of stochastic approximation algorithms for machine learning. In *NIPS*, pages 451–459, 2011.
- [6] Debraj Basu, Deepesh Data, Can Karakus, and Suhas N. Diggavi. Qsparse-local-sgd: Distributed SGD with quantization, sparsification, and local computations. *CoRR*, abs/1906.02367, 2019.
- [7] J. Bernstein, Y. Wang, K. Azizzadenesheli, and A. Anandkumar. SignSGD: compressed optimisation for non-convex problems. In *ICML*, pages 559–568, 2018.
- [8] L. Bottou. Large-scale machine learning with stochastic gradient descent. In COMPSTAT, pages 177–186, 2010.
- [9] Kai Chen and Qiang Huo. Scalable training of deep learning machines by incremental block training with intra-block parallel optimization and blockwise model-update filtering. In *ICASSP*, pages 5880–5884, 2016.
- [10] Gregory F. Coppola. *Iterative parameter mixing for distributed large-margin training of structured predictors for natural language processing*. PhD thesis, University of Edinburgh, UK, 2015.
- [11] R. Gitlin, J. Mazo, and M. Taylor. On the design of gradient algorithms for digitally implemented adaptive filters. *IEEE Transactions on Circuit Theory*, 20(2):125–136, March 1973.
- [12] Elad Hazan and Satyen Kale. Beyond the regret minimization barrier: optimal algorithms for stochastic strongly-convex optimization. *Journal of Machine Learning Research*, 15(1):2489– 2512, 2014.
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In CVPR, pages 770–778, 2016.
- [14] Robbins Herbert and Sutton Monro. A stochastic approximation method. *The Annals of Mathematical Statistics. JSTOR*, 22, no. 3:400–407, 1951.
- [15] Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian U. Stich, and Martin Jaggi. Error feedback fixes signsgd and other gradient compression schemes. In *ICML*, pages 3252–3261, 2019.
- [16] Anastasia Koloskova, Sebastian U. Stich, and Martin Jaggi. Decentralized stochastic optimization and gossip algorithms with compressed communication. In *ICML*, pages 3478–3487, 2019
- [17] Jakub Konecný. Stochastic, distributed and federated optimization for machine learning. CoRR, abs/1707.01155, 2017.

- [18] Maksim Lapin, Matthias Hein, and Bernt Schiele. Top-k multiclass SVM. In *NIPS*, pages 325–333, 2015.
- [19] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. In *Proceedings of the IEEE*, 86(11):2278-2324, 1998.
- [20] Y. Lin, S. Han, H. Mao, Y. Wang, and W. J. Dally. Deep gradient compression: Reducing the communication bandwidth for distributed training. In *ICLR*, 2018.
- [21] H. Mania, X. Pan, D. S. Papailiopoulos, B. Recht, K. Ramchandran, and M. I. Jordan. Perturbed iterate analysis for asynchronous stochastic optimization. *SIAM Journal on Optimization*, 27(4):2202–2229, 2017.
- [22] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas. Communication-efficient learning of deep networks from decentralized data. In *AISTATS*, pages 1273–1282, 2017.
- [23] Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. SIAM Journal on Optimization, 19(4):1574– 1609, 2009.
- [24] Lam M. Nguyen, Phuong Ha Nguyen, Marten van Dijk, Peter Richtárik, Katya Scheinberg, and Martin Takác. SGD and hogwild! convergence without the bounded gradients assumption. In *ICML*, pages 3747–3755, 2018.
- [25] A. Rakhlin, O. Shamir, and K. Sridharan. Making gradient descent optimal for strongly convex stochastic optimization. In *ICML*, 2012.
- [26] Benjamin Recht, Christopher Ré, Stephen J. Wright, and Feng Niu. Hogwild: A lock-free approach to parallelizing stochastic gradient descent. In NIPS, pages 693–701, 2011.
- [27] F. Seide, H. Fu, J. Droppo, G. Li, and D. Yu. 1-bit stochastic gradient descent and its application to data-parallel distributed training of speech dnns. In *INTERSPEECH*, pages 1058–1062, 2014.
- [28] A. Sergeev and M. D. Balso. Horovod: fast and easy distributed deep learning in tensorflow. *CoRR*, abs/1802.05799, 2018.
- [29] Shai Shalev-Shwartz, Yoram Singer, and Nathan Srebro. Pegasos: Primal estimated sub-gradient solver for SVM. In *ICML*, pages 807–814, 2007.
- [30] S. U. Stich, J. B. Cordonnier, and M. Jaggi. Sparsified SGD with memory. In *NeurIPS*, pages 4452–4463, 2018.
- [31] Sebastian U. Stich. Local SGD converges fast and communicates little. In ICLR, 2019.
- [32] Nikko Strom. Scalable distributed DNN training using commodity GPU cloud computing. In INTERSPEECH, pages 1488–1492, 2015.
- [33] A. Theertha Suresh, F. X. Yu, S. Kumar, and H. B. McMahan. Distributed mean estimation with limited communication. In *ICML*, pages 3329–3337, 2017.
- [34] H. Tang, S. Gan, C. Zhang, T. Zhang, and Ji Liu. Communication compression for decentralized training. In *NeurIPS*, pages 7663–7673, 2018.
- [35] Hanlin Tang, Shaoduo Gan, Ce Zhang, Tong Zhang, and Ji Liu. Communication compression for decentralized training. In *NeurIPS*, pages 7663–7673, 2018.
- [36] H. Wang, S. Sievert, S. Liu, Z. B. Charles, D. S. Papailiopoulos, and S. Wright. ATOMO: communication-efficient learning via atomic sparsification. In *NeurIPS*, pages 9872–9883, 2018.
- [37] Jianyu Wang and Gauri Joshi. Cooperative SGD: A unified framework for the design and analysis of communication-efficient SGD algorithms. *CoRR*, abs/1808.07576, 2018.
- [38] J. Wangni, J. Wang, J. Liu, and T. Zhang. Gradient sparsification for communication-efficient distributed optimization. In *NeurIPS*, pages 1306–1316, 2018.
- [39] W. Wen, C. Xu, F. Yan, C. Wu, Y. Wang, Y. Chen, and H. Li. Terngrad: Ternary gradients to reduce communication in distributed deep learning. In *NIPS*, pages 1508–1518, 2017.
- [40] J. Wu, W. Huang, J. Huang, and T. Zhang. Error compensated quantized SGD and its applications to large-scale distributed optimization. In *ICML*, pages 5321–5329, 2018.
- [41] Tianyu Wu, Kun Yuan, Qing Ling, Wotao Yin, and Ali H. Sayed. Decentralized consensus optimization with asynchrony and delays. *IEEE Trans. Signal and Information Processing over Networks*, 4(2):293–307, 2018.

- [42] Hao Yu, Rong Jin, and Sen Yang. On the linear speedup analysis of communication efficient momentum sgd for distributed non-convex optimization. In *ICML*, pages 7184–7193, 2019.
- [43] Hao Yu, Sen Yang, and Shenghuo Zhu. Parallel restarted SGD with faster convergence and less communication: Demystifying why model averaging works for deep learning. In *AAAI*, pages 5693–5700, 2019.
- [44] Y. Zhang, J. C. Duchi, M. I. Jordan, and M. J. Wainwright. Information-theoretic lower bounds for distributed statistical estimation with communication constraints. In NIPS, pages 2328–2336, 2013
- [45] Y. Zhang, J. C. Duchi, and M. J. Wainwright. Communication-efficient algorithms for statistical optimization. *Journal of Machine Learning Research*, 14(1):3321–3363, 2013.