# A Distributed Synchronous SGD Algorithm with Global Top-$k$ Sparsification for Low Bandwidth Networks

Shaohuai Shi[†], Qiang Wang[†], Kaiyong Zhao[†], Zhenheng Tang[†], Yuxin Wang[†], Xiang Huang[‡], Xiaowen Chu[*†]

[†]Department of Computer Science, Hong Kong Baptist University

[‡]MassGrid.com, Shenzhen District Block Technology Co., Ltd.

{csshshi, qiangwang, kyzhao, zhtang, yxwang}@comp.hkbu.edu.hk, xiang.huang@galasports.net, chxw@comp.hkbu.edu.hk

*Abstract*—Distributed synchronous stochastic gradient descent (S-SGD) with data parallelism has been widely used in training large-scale deep neural networks (DNNs), but it typically requires very high communication bandwidth between computational workers (e.g., GPUs) to exchange gradients iteratively. Recently, Top-$k$ sparsification techniques have been proposed to reduce the volume of data to be exchanged among workers and thus alleviate the network pressure. Top-$k$ sparsification can zero-out a significant portion of gradients without impacting the model convergence. However, the sparse gradients should be transferred with their indices, and the irregular indices make the sparse gradients aggregation difficult. Current methods that use All-Gather to accumulate the sparse gradients have a communication complexity of $O(kP)$, where $P$ is the number of workers, which is inefficient on low bandwidth networks with a large number of workers. We observe that not all top-$k$ gradients from $P$ workers are needed for the model update, and therefore we propose a novel global Top-$k$ (gTop-$k$) sparsification mechanism to address the difficulty of aggregating sparse gradients. Specifically, we choose global top-$k$ largest absolute values of gradients from $P$ workers, instead of accumulating all local top-$k$ gradients to update the model in each iteration. The gradient aggregation method based on gTop-$k$ sparsification, namely gTopKAllReduce, reduces the communication complexity from $O(kP)$ to $O(k \log P)$. Through extensive experiments on different DNNs, we verify that gTop-$k$ S-SGD has nearly consistent convergence performance with S-SGD, and it has only slight degradations on generalization performance. In terms of scaling efficiency, we evaluate gTop-$k$ on a cluster with 32 GPU machines which are interconnected with 1 Gbps Ethernet. The experimental results show that our method achieves $2.7 - 12\times$ higher scaling efficiency than S-SGD with dense gradients and $1.1 - 1.7\times$ improvement than the existing Top-$k$ S-SGD.

*Index Terms*—Deep Learning; Stochastic Gradient Descent; Distributed SGD; Gradient Communication; Top-$k$ Sparsification; gTop-$k$;

## I. INTRODUCTION

With the increase of training data volume and the growing complexity of deep neural networks (DNNs), distributed computing environments (such as GPU clusters) are widely adopted to accelerate the training of DNNs. The data-parallel synchronous stochastic gradient descent (S-SGD) method is one of the commonly used optimizers to minimize the objective function of large-scale DNNs [1][2]. Compared to SGD on a single worker, S-SGD distributes the workloads

*Corresponding author.

TABLE I
COMMUNICATION COMPLEXITY OF GRADIENT AGGREGATION ALGORITHMS

| Aggregation Algorithm | Complexity | Time Cost |
|---|---|---|
| DenseAllReduce | $O(m)$ | $2(P-1)\alpha + 2\frac{P-1}{P}m\beta$ |
| TopKAllReduce | $O(kP)$ | $\log(P)\alpha + 2(P-1)k\beta$ |
| Ours (gTopKAllReduce) | $O(k \log P)$ | $2\log(P)\alpha + 4k\log(P)\beta$ |

Note: $m$ is the number of model parameters. $P$ is the number of workers. $k = \rho \times m$ is the number of local gradients to be aggregated. $\alpha$ and $\beta$ are system dependent constants.

to multiple workers to accelerate the training, but it also introduces the communication overhead of exchanging model parameters or gradients in each iteration. Assume that there are $P$ workers training a single DNN model with S-SGD. In every iteration, all workers take different mini-batches of data to calculate the model gradients in parallel. Then they need to average the gradients before updating the model parameters, which involves significant data communications [3]. Since the computing power of computational units (e.g., GPUs and TPUs) grows much faster than the growth of network speed, network communication performance has now become the training bottleneck, especially when the communication-to-computation ratio is high [4]. Many large IT companies use expensive high-speed networks such as 40/100Gbps IB or Ethernet to alleviate the communication pressure, but still many researchers and small companies can only use consumer-level GPUs connected by low-bandwidth networks such as 1Gig-Ethernet.

To conquer the communication challenge, one can either increase the workload of workers by choosing a large batch size or reduce the required data communications in each iteration. Very recently, many large-batch SGD techniques have been proposed with sophisticated optimization strategies [5][6][7][8][9] to increase the scaling efficiency without losing the model accuracy. On the other hand, gradient sparsification, quantification and compression methods [10][11][12][13][14][15][16] have been proposed to dramatically reduce the size of exchanged gradients without affecting the convergence rate. Among the model/gradient size reduction techniques, the Top-$k$ sparsification is one of the key approaches [17][12][18] that can sparsify the local gradients

to just about $0.001$ density ($99.9\%$ gradients are zeros and there is no need to transfer these zero-out values) [11][12].

Top-$k$ sparsification has been a successful gradient compression method with empirical and theoretical studies in [17][12][16], in which researchers have verified that only a small number of gradients are needed to be averaged during the phase of gradient aggregation without impairing model convergence or accuracy. However, the sparsified gradients are generally associated with irregular indices, which makes it a challenge to accumulate the selected gradients from all workers[1] efficiently. The ring-based AllReduce method used on dense gradients (DenseAllReduce) has an $O(P + m)$ communication complexity [19], where $P$ is the number of workers and $m$ is the size of parameters (or gradients). In Top-$k$ sparsification, assume that the density of gradients is $\rho$ on each worker, $k = \rho \times m$, and the corresponding indices of non-zero values are irregular from different workers and iterations; thus it generally needs to transfer $2k$ number of values (gradients and indices) in each iteration. However, with the sparse gradients, the DenseAllReduce method cannot be directly used to accumulate all the sparse gradients with irregular indices, and a recent solution uses the AllGather collective [20], which requires an $O(kP)$ communication complexity. We use TopKAllReduce to denote the method of averaging irregularly indexed top-$k$ gradients by adopting AllGather. When scaling to a large number of workers (i.e., $P$ is large), even high sparsification ratios still generate significant communication overheads.

The main idea of Top-$k$ sparsification is based on the fact that gradients with larger absolute values can contribute more to the model convergence. Theoretical analysis on this idea has been proposed in [21][16][18]. Therefore, one can further select Top-$k$ gradients from the accumulated results from $P$ groups of top-$k$ values generated by $P$ workers. That is, even though $P$ workers can generate a maximum number of $k \times P$ non-zero gradients after aggregation, the top-$k$ gradients (in terms of absolute values) would be the most important for the model updates. Based on this observation, we propose an efficient Top-$k$ sparsification method to tackle the difficulty of TopKAllReduce with little impacting on the model convergence. Specifically, instead of accumulating the irregularly indexed non-zero gradients from all workers, we choose the global Top-$k$ (gTop-$k$) gradients in terms of absolute values. gTop-$k$[2] can elegantly make use of a tree structure to select the global top-$k$ values from all workers, which we call gTopKAllReduce, such that the communication complexity is reduced from $O(kP)$ to $O(k \log P)$. We summarize the communication complexities of different gradient aggregation solutions in Table I.

In this paper, we first implement the gTopKAllReduce algorithm which provides much more efficient global Top-$k$ sparse gradients aggregation from distributed workers. Then

we integrate our proposed gTopKAllReduce to gTop-$k$ S-SGD under PyTorch[3], which is one of the most popular deep learning frameworks and MPI[4]. On a 32-node GPU cluster connected by 1-Gbps Ethernet, gTop-$k$ S-SGD achieves 2.7-12.8x speedup than S-SGD with highly optimized libraries Horovod [22] and NCCL[5]. Compared to Top-$k$ S-SGD, gTop-$k$ S-SGD is generally around $1.5$ times faster on the evaluated experiments on various DNNs and datasets. Our contributions are summaries as follows:

- We observe that the accumulating results of Top-$k$ sparsification can be further sparsified before being updated to the model.
- We propose an efficient global Top-$k$ sparsification algorithm on distributed SGD, called gTop-$k$ S-SGD, to accelerate distributed training of deep neural networks.
- We implement the proposed gTop-$k$ S-SGD atop popular framework PyTorch and MPI, and we also release all our experimental parameters for reproducibility[6].
- gTop-$k$ S-SGD achieves significantly improved scaling efficiency on the real-world applications with various DNNs and datasets under low-bandwidth networks (e.g., 1 Gig-Ethernet).

The rest of the paper is organized as follows. We introduce the preliminaries in Section II, in which some background information and the main problem is clarified. In Section III, we present our observation from Top-$k$ sparsification and propose an efficient gTop-$k$ S-SGD algorithm. Then we demonstrate the detailed experimental study in Section IV and have a discussion in Section V. Section VI gives an introduction to the related work, and finally we conclude the paper in Section VII.

## II. PRELIMINARIES

### A. DNNs

Deep neural networks (DNNs) are generally stacked with many hierarchical layers, and each layer is a transformer function of the input values. We can formulate an $L$-layer DNN by

$$a^{(l)} = f(W^{(l)}, x^{(l)}), \qquad (1)$$

where $x^{(l)}$ and $a^{(l)}$ are the input and output of layer $l$ ($l = 1, 2, ..., L$) respectively. Inputs of current layer are the outputs of its previous layer(s) (e.g., $x^l = a^{(l-1)}$). The function $f$ is the transformer function which consists of an operation (e.g., inner product or convolution) and an activation function (e.g., ReLU). $W^{(l)}$ are the trainable model parameters, which could be iteratively updated during the training process using mini-batch stochastic gradient descent (SGD) optimizers and the backpropagation algorithm.

---

[1]Worker and GPU are interchangeable in this paper.

[2]In this paper, we mainly discuss the decentralized S-SGD with AllReduce to apply gTop-$k$ sparsification. But it is also applicable to the Parameter Server based distributed SGD.

[3]https://pytorch.org/

[4]https://www.open-mpi.org/

[5]https://developer.nvidia.com/nccl

[6]All experimental settings and source codes can be found at GitHub: https://github.com/hclhkbu/gtopkssgd

## B. Mini-batch SGD

The objective function $\mathcal{L}(W, D)$ defines the differences between the prediction values by the DNN and the ground truth. The mini-batch SGD optimizer updates the parameters iteratively to minimize the objective function. To be specific, there are three phases in each iteration during training: 1) Feed-forward phase: a mini-batch of data $D_i$ ($D_i \subset D$) is read as inputs of a DNN, and $D_i$ is fed forward across the neural network from the first layer to the last layer, which finally generates the prediction values to be used by the objective function $\mathcal{L}(W, D)$. 2) Backward-propagation phase: the gradients w.r.t. the parameters and inputs are calculated from the last layer to the first layer. 3) Update phase, the parameters are updated by the afore-generated gradients using the following formula (or its variants):

$$W_{i+1} = W_i - \eta \cdot \nabla \mathcal{L}(W_i, D_i), \qquad (2)$$

where $\eta$ is the learning rate. For a single-worker training, phases 1) and 2) are the main time costs of an iteration, which are computing-intensive tasks. So the average time of one iteration can be approximated by $t_{iter} = t_f + t_b$.

## C. Synchronous SGD

Synchronous SGD (S-SGD) with data parallelism is widely applied to train models with multiple workers (say $P$ workers, and indexed by $g$). Each worker keeps a consistent model at the beginning of each iteration. The workers take different mini-batches of data $D_i^g$ and forward them by phase 2), and then follow phase 3) to calculate the gradients $\nabla \mathcal{L}(W_i, D_i^g)$ in parallel. The average gradients from different workers are used to update the model. The update formula of parameters is rewritten as

$$W_{i+1} = W_i - \eta \frac{1}{P} \sum_{g=1}^{P} G_i^g, \qquad (3)$$

where $G_i^g = \nabla \mathcal{L}(W_i, D_i^g)$ denotes the gradients of worker $g$ at the $i^{th}$ iteration. The gradients are located in different workers without shared memory so that the averaging operation of gradients involves communication costs, which could become another system bottleneck. The average iteration time of S-SGD can be approximated by $t_{iter} = t_f + t_b + t_c$. Assume that we use weak-scaling on $P$ workers with S-SGD, the scaling efficiency can be approximated by

$$e = \frac{t_f + t_b}{t_f + t_b + t_c}. \qquad (4)$$

$t_c$ is generally related to $P$ and the model/gradient size $m$.

## D. DenseAllReduce

In Eq. 3, the summation of $G_i^g$ (i.e., $\sum_{g=1}^{P} G_i^g$) can be directly implemented by an AllReduce collective, which is denoted as DenseAllReduce. The ring-based AllReduce algorithm (which is also included in NCCL) is an efficient implementation on the dense-GPU cluster. To understand the time cost of DenseAllReduce, we revisit the time model of the ring-based AllReduce. According to [23], the time cost of ring-based AllReduce can be represented by

$$t_c^{dar} = 2(P-1)\alpha + 2\frac{P-1}{P}m\beta, \qquad (5)$$

where $\alpha$ is the latency (startup time) of a message transfer between two nodes, and $\beta$ is the transmission time per element between two nodes using the alpha-beta model [24].

---

**Algorithm 1** S-SGD with Top-$k$ sparsification on worker $g$ [12][20]

---

**Input:** The dataset: $D$
The initialized weights: $W$
The mini-batch size per worker: $b$
The number of workers: $P$
The number of iterations to train: $N$
The number gradients to select: $k$

1: $G_0^g = 0$
2: **for** $i = 1 \rightarrow N$ **do**
3:     Sampling a mini-batch of data $D_i^g$ from $D$;
4:     $G_i^g = G_{i-1}^g + \nabla \mathcal{L}(W_i, D_i^g)$;
5:     Select threshold $thr =$ the $k^{th}$ largest value of $|G_i^g|$;
6:     $Mask = |G_i^g| > thr$;
7:     $\widetilde{G}_i^g = G_i^g \odot Mask$; // Mask has $k$ non-zero values
8:     $G_i^k = G_i^g \odot \neg Mask$; // Store the residuals
9:     $G_i =$TopKAllReduce($\widetilde{G}_i^g$); //$G_i = \frac{1}{P}\sum_{g=1}^{P} \widetilde{G}_i^g$
10:     $W_i = W_{i-1} - \eta G_i$;
11: **end for**
12: **procedure** TopKAllReduce($\widetilde{G}_i^g$)
13:     $[\mathcal{V}_i^g, \mathcal{I}_i^g] = \widetilde{G}_i^g$;
14:     $[\mathcal{V}, \mathcal{I}] =$AllGather($[\mathcal{V}_i^g, \mathcal{I}_i^g]$);
15:     $G_i = \widetilde{G}_i^g$;
16:     **for** $g = 0 \rightarrow P - 1$ **do**
17:         $G_i[\mathcal{I}[g * P : g * (P+1)]] + = V[g * P : g * (P+1)]$;
18:     **end for**
19:     $G_i = G_i / P$;
20:     Return $G_i$;
21: **end procedure**

---

## E. Top-k sparsification

From Eq. 5, it is noted that with $P$ or $m$ becoming large, the communication cost will be linearly increased. To reduce the size of transfer messages $m$, researchers propose Top-$k$ sparsification [12] which introduces highly sparse gradients. With Top-$k$ sparsification, each worker only needs to contribute the $k$ largest absolute values of gradients $G_i^g$ to be summed up in each iteration, and the zeroed-out values of gradients are stored locally and accumulated at the next iteration. Both theoretical and empirical studies have verified that the Top-$k$ sparsification has little impact on the model convergence and accuracy [17][12][16]. For completeness, the pseudo-code of Top-$k$ sparsification S-SGD is shown in Algorithm 1. Note that at Line 9 of Algorithm 1, the implementation of TopKAllReduce is completely different from the DenseAllReduce for efficiency since the non-zero values of $\widetilde{G}_i^g$ may come from inconsistent indices $\mathcal{I}_i^g$ from different workers. Efficient implementations of such sparse AllReduce are non-trivial. Current methods [20] are using AllGather to implement TopKAllReduce, in which the sparsified gradients

are gathered as a dense vector combined with its corresponding indices, say $\widetilde{G}_i^g = [\mathcal{V}_i^g, \mathcal{I}_i^g]$. Both sizes of $\mathcal{V}_i^g$ and $\mathcal{I}_i^g$ are $k$. According to the communication model of AllGather [19], the time cost for all-gathering $2k$ values is

$$t_c^{tar} = \log(P)\alpha + 2(P-1)k\beta. \qquad (6)$$

From Eq. 6, we can see that with increasing $P$, $t_c^{tar}$ is linearly increased. Therefore, the effect of Top-$k$ sparsification will diminish with the increase of number of workers. In this paper, we propose a global Top-$k$ (gTop-$k$) sparsification algorithm to address this scalability problem.

## III. METHODOLOGY

In this section, we first demonstrate some observations from Top-$k$ sparsification S-SGD, and then we present our proposed global Top-$k$ sparsification algorithm. For ease of presentation, we assume that the number of workers $P$ is the power of 2.

### A. Observations from Top-k sparsification

In the previous section, we have introduced Top-$k$ sparsification S-SGD, in which there are $k$ values selected from the local worker and then are accumulated across all the workers. We get insight into the distribution of non-zero values (denoted as $\mathcal{G}_i$) of $G_i$ which is generated by the summation of the sparse gradients from all workers. We found that not all $\mathcal{G}_i$ (whose number of elements is $\mathcal{K}$, and $k \leq \mathcal{K} \leq k \times P$) contribute to the model convergence. Specifically, $\mathcal{G}_i$ can be further sparsified as $\widetilde{\mathcal{G}}_i$ such that only a smaller number of non-zero gradients are needed for model updates. In other words, one can further select top-$k$ largest absolute values, $\widetilde{\mathcal{G}}_i$, from $\mathcal{G}_i$ to update the model while maintaining the model convergence. In this scenario, the selected $\widetilde{\mathcal{G}}_i$ from $\mathcal{G}_i$ results in the fact that $\mathcal{K}-k$ afore-summed gradients are neither updated to the model nor stored into the local residuals, which could damage the model convergence. Therefore, if only $k$ elements are selected from $\mathcal{G}_i$ to update the model, the remaining $\mathcal{K}-k$ elements should be put back as residuals with corresponding indices so that they can be accumulated locally and eventually contribute to the model update in some future iterations. We empirically verify this idea by training a ResNet DNN, and show the convergence result in Fig. 1.
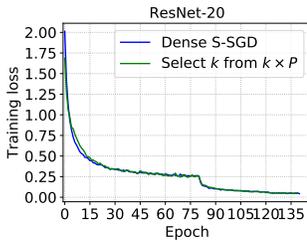


Fig. 1. The convergence of ResNet-20 on 4 workers with only $k = 0.001 \times m$ elements updated at each iteration. The experimental settings can be found in Section IV.

---

**Algorithm 2** Naive version S-SGD with gTop-$k$ on worker $g$

**Input:** The dataset: $D$
The initialized weights: $W$
The mini-batch size per worker: $b$
The number of workers: $P$
The number of iterations to train: $N$
The number of gradients to select: $k$

1: $G_0^g = 0$
2: **for** $i = 1 \to N$ **do**
3:     Sampling a mini-batch of data $D_i^g$ from $D$;
4:     $G_i^g = G_{i-1}^g + \nabla\mathcal{L}(W_i, D_i^g)$;
5:     Select threshold $thr$ = the $k^{th}$ largest value of $|G_i^g|$;
6:     $Mask = |G_i^g| > thr$;
7:     $\widetilde{G}_i^g = G_i^g \odot Mask$; // Mask has $k$ non-zero values
8:     $G_i^g = G_i^g \odot \neg Mask$; // Store the residuals
9:     $G_i$ =SparseAllReduce($\widetilde{G}_i^g$); //$G_i = \frac{1}{P}\sum_{g=1}^{P} \widetilde{G}_i^g$
10:     // At this time all workers have consistent $G_i$
11:     Select global threshold $gThr$ = the $k^{th}$ largest value of $|G_i|$;
12:     $gMask = |G_i| > gThr$;
13:     $\widetilde{G}_i = G_i \odot gMask$;
14:     $G_i^g += \widetilde{G}_i^g \odot \neg gMask \odot Mask$; // Store extra residuals
15:     $W_i = W_{i-1} - \eta\widetilde{G}_i$;
16: **end for**

---

### B. The key idea of gTop-k

According to the above observations, it only needs $k$ largest absolute values from all the sparsified gradients $G_i^g$, where $g = 1, 2, ..., P$. Therefore, the problem is formulated as the global Top-$k$ (gTop-$k$) selection from $G_i$ instead of using TopKAllReduce, while $\widetilde{G}_i^g$ are located in distributed workers. We again let $[V_i^g, \mathcal{I}_i^g]$ denote the non-zero values and corresponding indices of $\widetilde{G}_i^g$ whose number of non-zero values is $k$. We first use the AllGather version to illustrate the key idea of gTop-$k$ sparsification, and then we present our new efficient algorithm for gTop-$k$ sparsification. At Line 10 of Algorithm 1, $W_i = W_{i-1} - \eta G_i$, $G_i$ with $\mathcal{K}$ non-zero values contributing updates to $W_i$. We further sparsify $G_i$ by selecting $k$ largest absolute values from $G_i$. The implementation is shown in Algorithm 2. Please be noted that this version is only used to illustrate the key idea of how to select those gradients to update the model. The efficient algorithm is presented afterwards in the next subsection. An example of gTop-$k$ sparsification using AllGather on 4 workers is shown in Fig. 2.
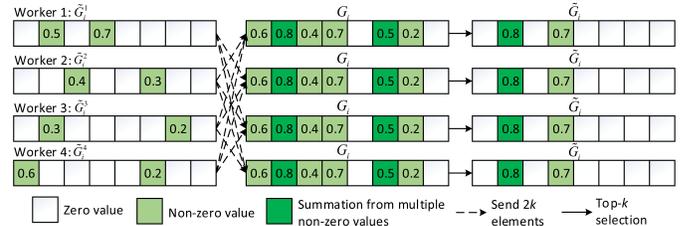


Fig. 2. An example of gTop-$k$ using AllGather on 4 workers, and $k = 2$.

## C. gTopKAllReduce: An efficient AllReduce algorithm for gTop-k sparsification

From Eq. 6, we can see that the AllGather collective is inefficient to conduct the AllReduce operation from irregular indexed gradients. Based on the same density, the main purpose of our proposed efficient algorithm is to alleviate the high impact of the variable $P$ on the time cost. For ease of presentation, we first define a Top-$k$ operation, $\top$, of two sparse vectors, say $\widetilde{G}^a$ and $\widetilde{G}^b$, both of which have $k$ non-zero values.

*Definition 1:* A Top-$k$ operation: $\top$. $\widetilde{G}^{a,b} = \widetilde{G}^a \top \widetilde{G}^b = mask \odot (\widetilde{G}^a + \widetilde{G}^b)$, where $mask = (|\widetilde{G}^a + \widetilde{G}^b| > thr)$, and $thr =$ the $k^{th}$ largest value of $|\widetilde{G}^a + \widetilde{G}^b|$.

Note that the number of non-zero values of $\widetilde{G}^{a,b}$ is also $k$. During the training of S-SGD, $\widetilde{G}^a$ and $\widetilde{G}^b$ are located in different workers without shared memory. One should exchange the two sparse vectors to achieve a global Top-$k$ sparse vector: $\widetilde{G}^{a,b}$. The operation for two distributed workers is shown in Fig. 3, which demonstrates that $\top$ can be efficiently implemented by a send operation (network communication), followed by a local Top-$k$ selection on a maximum number of $2k$ non-zero values.
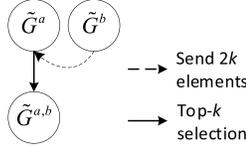


Fig. 3. An implementation of $\top$ for two distributed sparse vectors $\widetilde{G}^a$ and $\widetilde{G}^b$. The second worker ($[V^b, \mathcal{I}^b] = \widetilde{G}^b$) with $k$ non-zero elements ($V^b$) combined with $k$ indices ($\mathcal{I}^b$) sends $[V^b, \mathcal{I}^b]$ to the first worker, and then the first worker has the information of indices to add the values received from the second worker, i.e., $\widetilde{G}^a + \widetilde{G}^b$, and the first worker easily computes $\widetilde{G}^{a,b} = \widetilde{G}^a \top \widetilde{G}^b$ according to Definition 1.

When scaling to $P$ workers (assume that $P$ is the power of 2), since the final $k$ is equal to the local $k$, we propose a tree structure based technique to reduce the total transfer size. To illustrate the tree structure for gTop-$k$, we show an 8-worker example in selecting the global Top-$k$ values in Fig. 4. There are 3 rounds of communications for 8 workers (i.e., $\log_2 8 = 3$). At the $j^{th}$ round, there are $\frac{P}{2^j}$ pairs of workers to do the $\top$ operations in parallel. After 3 rounds, the first worker (rank 0) finally generates the global Top-$k$ values (i.e., $\widetilde{G} = \widetilde{G}^{1,2,...,8} = \widetilde{G}^1 \top \widetilde{G}^2 \top ... \widetilde{G}^8$).

According to the illustration of tree structure based gTop-$k$, we propose the gradients aggregation with gTop-$k$ sparsification, which is called gTopKAllReduce in Algorithm 3. Line 1 selects the non-zero values from sparse $\widetilde{G}^g$ to assign the variable "$sends$", which should be sent to other workers. Line 2 allocates the buffer "$recvs$" to receive the "$sends$" from another worker at each communication round. Lines 4-17 describe the procedure of $\widetilde{G} = \widetilde{G}^1 \top \widetilde{G}^2 \top ... \widetilde{G}^P$, which is finally stored in "$sends$" for the next round communication. The functions "Recv" and "Send" in Line 10 and 14 are a paired operation and can be implemented by MPI. Since the
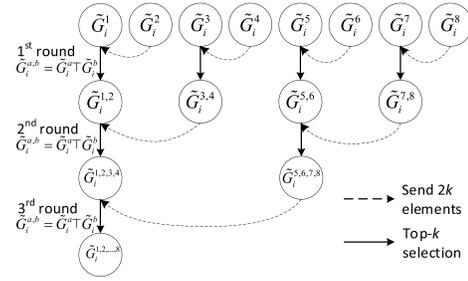


Fig. 4. An example of gTop-$k$ for 8 distributed sparse vectors $\widetilde{G}^1, \widetilde{G}^2, ..., \widetilde{G}^8$. I.e., $\widetilde{G} = \widetilde{G}^{1,2,...,8} = \widetilde{G}^1 \top \widetilde{G}^2 \top ... \widetilde{G}^8$. It only requires $log_2 P = log_2 8 = 3$ rounds of network communications to select the global Top-$k$.

---

**Algorithm 3** gTopKAllReduce

---

**Input:** The sparsified gradients: $\widetilde{G}^g$
The number of non-zero elements: $k$
The number of workers: $P$
The rank of worker: $g$
1: $sends = [V^g, \mathcal{I}^g] = \widetilde{G}^g$;
2: Initialize $recvs$ with the same as $sends$;
3: $nRounds = \log P$;
4: **for** $i = 1 \rightarrow nRounds$ **do**
5:     $participateRanks = [1 \rightarrow P, step = i]$;
6:     **if** $g$ in $participateRanks$ **then**
7:         $localRank = participateRanks.index(g)$;
8:         **if** $localRank\%2 == 0$ **then**
9:             $source = participateRanks[localRank + 1]$;
10:            Recv($recvs$, source=$source$);
11:            $sends = recvs \top sends$;
12:         **else**
13:            $target = participateRanks[localRank - 1]$;
14:            Send($sends$, dest=$target$);
15:         **end if**
16:     **end if**
17:     Barriar();
18: **end for**
19: Bcast($sends$, root=0);
20: $\widetilde{G} = [V, \mathcal{I}] = sends$;
21: $Mask = [0]$ in the same number of elements with $\widetilde{G}$;
22: $Mask[\mathcal{I}] = 1$;
23: Return $\widetilde{G}, Mask$;

---

result $\widetilde{G}$ by far is only stored at the first worker (rank=0), Line 19 broadcasts the $\widetilde{G}$ to all other workers, which also requires $\log P$ number of communications using the flat-tree algorithm [25]. Finally, Lines 21 and 22 record the $Mask$ which indicates the indices that are used in $\widetilde{G}$.

### D. Communication complexity analysis of gTopKAllReduce

There are two main processes of gTopKAllReduce. The first one is the calculation of $\widetilde{G}$. From Fig. 4, the first worker should take part in the communication at every round, so we only need to analyze the communication complexity of the worker with rank 0. Rank 0 takes $\log P$ rounds of communications to calculate $\widetilde{G}$, and it receives $2k$ elements from another worker at each round which takes a time cost of $\alpha + 2k\beta$. Thus, the overall time cost of the first process is $\alpha \log P + 2k\beta \log P$. In the second process, the global top-$k$ values (i.e., $\widetilde{G}$) in the first worker should be broadcasted to all the other workers.

The broadcast operation takes $\alpha \log P + 2k\beta \log P$ according to the flat-tree algorithm. In summary, the time cost of gTopKAllReduce is

$$t_c^{gar} = 2\alpha \log P + 4k\beta \log P. \qquad (7)$$

The communication complexity is much lower than Top-KAllReduce especially when $P$ is large.

*E. gTop-k S-SGD with gTopKAllReduce*

With the above proposed efficient implementation of gTop-KAllReduce, we improve the gTop-$k$ S-SGD in Algorithm 2 by replacing Lines 9-13 with a line that invokes gTopKAllReduce shown in Algorithm 3. The improved version of the gTop-$k$ S-SGD training algorithm is presented in Algorithm 4. Compared to Top-$k$ S-SGD, gTop-$k$ S-SGD only introduces an extra computation (Line 10 in Algorithm 4) whose overhead is much smaller than the communication overhead, while gTop-$k$ S-SGD reduces the communication complexity a lot.

---

**Algorithm 4** gTopKAllReduce based S-SGD on worker $g$

---

**Input:** The dataset: $D$
The initialized weights: $W$
The mini-batch size per worker: $b$
The number of workers: $P$
The number of iterations to train: $N$
The number gradients to select: $k$
1: $G_0^g = 0$
2: **for** $i = 1 \to N$ **do**
3:     Sampling a mini-batch of data $D_i^g$ from $D$;
4:     $G_i^g = G_{i-1}^g + \nabla\mathcal{L}(W_i, D_i^g)$;
5:     Select threshold $thr =$ the $k^{th}$ largest value of $|G_i^g|$;
6:     $Mask = |G_i^g| > thr$;
7:     $\widetilde{G}_i^g = G_i^g \odot Mask$; // Mask has $k$ non-zero values
8:     $G_i^g = G_i^g \odot \neg Mask$; // Store the residuals
9:     $\widetilde{G}_i, gMask =$gTopKAllReduce($\widetilde{G}_i^g, k, P, g$);
10:     $G_i^g += \widetilde{G}_i^g \odot \neg gMask \odot Mask$; // Store extra residuals
11:     $W_i = W_{i-1} - \eta\widetilde{G}_i$;
12: **end for**

---

## IV. EXPERIMENTAL STUDY

We conduct extensive experiments to evaluate the effectiveness of our proposed gTop-$k$ S-SGD by real-world applications on a 32-GPU cluster. We first validate the convergence of gTop-$k$ S-SGD. Then we evaluate the time cost and efficiency of gTopKAllReduce and compare them with the dense AllReduce (DenseAllReduce) and Top-$k$ AllReduce (gTopKAllReduce) counterparts. After that, we make a comparison on the training efficiency among the three S-SGD algorithms (i.e., S-SGD with dense gradients, Top-$k$ S-SGD, and gTop-$k$ S-SGD). We also break down the training process in an iteration to several phases to analyze the extra overhead introduced by gTop-$k$ sparsification.

*A. Experimental setup*

**Hardware**: The distributed environments are configured as a 32-node cluster, each with one Nvidia P102-100 GPU. All nodes are connected by a 1-Gbps Ethernet. Details of the hardware are shown in Table II.

**Software**: All GPU machines are installed with Ubuntu-16.04, Nvidia GPU driver at version 390.48, and CUDA-9.1. The communication libraries are OpenMPI-3.1.1[7] and NCCL-2.1.5[8]. We use the highly optimized distributed training library Horovod[9] [22] at version 1.4.1. The deep learning framework is PyTorch at version 0.4.0 with cuDNN-7.1.

TABLE III
DEEP MODELS FOR TRAINING.

| Model | Dataset | # of Epochs | $b$ | $\eta$ |
|---|---|---|---|---|
| VGG-16 | Cifar-10 | 140 | 128 | 0.1 |
| ResNet-20 | Cifar-10 | 140 | 128 | 0.1 |
| AlexNet | ImageNet | 45 | 64 | 0.01 |
| ResNet-50 | ImageNet | 45 | 256 | 0.01 |
| LSTM-PTB | PTB | 40 | 100 | 1.0 |

Note: All models are trained with 32-bit floating points.

**DNNs**: We choose various DNNs from several areas of AI applications with different data sets. The data sets include Cifar-10 [26] that contains $50,000$ training samples, ImageNet [27] that contains about 1.2 million samples for image classification, and the Penn Treebank corpus (PTB) [28] that contains $923,000$ training samples for language modeling. For the Cifar-10 data set, we use the VGG-16 model [29] and the ResNet-20 model [30]. For the ImageNet data set, the AlexNet model [31] and the ResNet-50 model [30] are used. We exploit a 2-layer LSTM language model (LSTM-PTB) for the PTB dataset, which is similar as in [12]. The details of the deep learning models are given in Table III. We use momentum SGD with a momentum of 0.9 to train all models.

*B. Convergence comparison*

The convergence of Top-$k$ sparsification S-SGD has been verified to be nearly consistent with the dense version in much previous work [17][12][16], so we would not include the convergence curves of Top-$k$ S-SGD. We compare our gTop-$k$ S-SGD with the original S-SGD with dense gradients running on 4 workers. It has been shown that the warmup strategy in the first several epochs helps the model converge better [12], so we adopt a similar warmup configuration. To be specific, the first 4 epochs use the dynamic densities of $[0.25, 0.0725, 0.015, 0.004]$ and small learning rates, and the remaining epochs adopt a density of 0.001 (for CNNs) or 0.005 (for LSTM).

---

[7]https://www.open-mpi.org/
[8]https://developer.nvidia.com/nccl
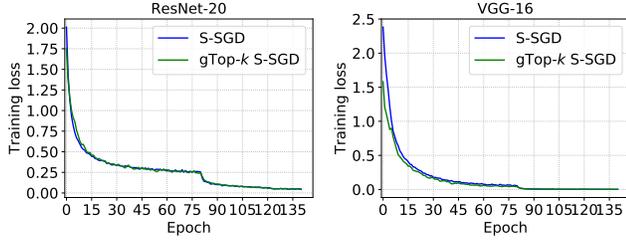[9]https://github.com/uber/horovod

Fig. 5. The convergences of VGG-16 and ResNet-20 with $P = 4$.

**Convergence on the Cifar-10 data set**: The convergences of VGG-16 and ResNet-20 models are shown in Fig. 5. The results show that the convergence rate of ResNet-20 is almost the same as the baseline, while the VGG-16 model even converges slightly better than the baseline.
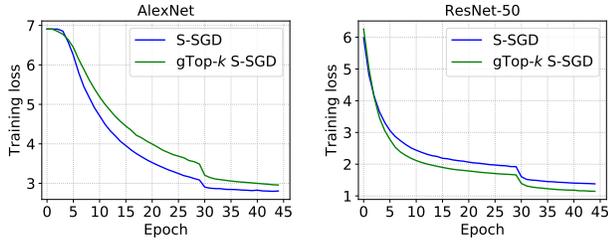


Fig. 6. The convergences of AlexNet and ResNet-50 with $P = 4$.

**Convergence on the ImageNet data set**: The convergences of AlexNet and ResNet-50 models are shown in Fig. 6. The results show that the convergence rates of the two CNNs are close to the baselines. On the AlexNet model, the convergence of gTop-$k$ S-SGD with $\rho = 0.001$ is slightly worse than the baseline, which could be caused by unbalanced parameters between convolutional layers and fully connected layers with the same low density. On the other hand, gTop-$k$ sparsification works very well on the ResNet-50 model, which converges even faster than the baseline.

**Convergence on the LSTM network**: The convergence of LSTM-PTB on the PTB data set is shown in Fig. 7. It is noted that the convergence of gTop-$k$ S-SGD is almost the same as that of S-SGD under a density of $0.005$.
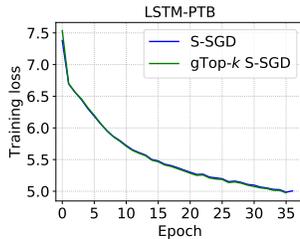


Fig. 7. The convergence of LSTM with $P = 4$ and $\rho = 0.005$.

In summary, three different types of DNNs from different data sets show that our proposed gTop-$k$ S-SGD would not

damage the model during training and keeps very close model convergence to dense S-SGD.

### C. Communication speed

To set the baseline, we first test the point-to-point communication performance with various sizes of messages because the performance of point-to-point communication plays an essential role in MPI collectives. Then we compare the communication performance of TopKAllReduce and gTopKAllReduce in different sizes of sparse vectors according to Table I.

**Point-to-point communication**: We test the point-to-point communication speed by using OSU Micro-Benchmark[10] at the version $5.5$. The time costs of the point-to-point communication between two machines are shown in Fig. 8, in which we run $5$ times to calculate the mean and standard variance from the reported values. It can be seen that the time used for transferring a message is a linear function with the size of the message, which verifies the $\alpha$-$\beta$ model. Based on the measured data, we can get $\alpha = 0.436ms$ and $\beta = 3.6 \times 10^{-5}ms$.
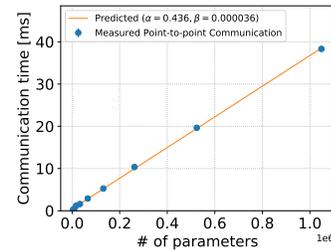


Fig. 8. Data transfer time in milliseconds with respective to the size of message on our experiment cluster.

**Time performance of AllReduce operations**: Since $P$ and $m$ are two main factors affecting the performance of TopKAllReduce and gTopKAllReduce, we compare their time performances in two dimensions (i.e., $P$ and $m$) based on the measured $\alpha$, $\beta$ and the time cost models in Table. I. First, we compare the time cost with different number of workers (from 4 to 128) based on $m = 25 \times 10^6$ (the approximate model size of ResNet-50) and $\rho = 0.001$. Second, in the configuration of a cluster with 32 workers, we make a comparison on how the time cost changes with the size of parameters increases. The time comparison is shown in Fig. 9. From the left of Fig. 9, when the number of nodes is small, TopKAllReduce is slightly faster than gTopKAllReduce. However, when the number of nodes increases to 16, TopKAllReduce becomes much worse than gTopKAllReduce. Furthermore, our proposed gTopKAllReduce is much more efficient than TopKAllReduce when scaling to large sizes of messages. To summarize, a larger number of workers or a larger message size would make gTopKAllReduce more efficient than TopKAllReduce.

### D. Scaling efficiency

The scaling efficiency of S-SGD with three different AllReduce algorithms are shown in Fig. 10. It can be seen that the

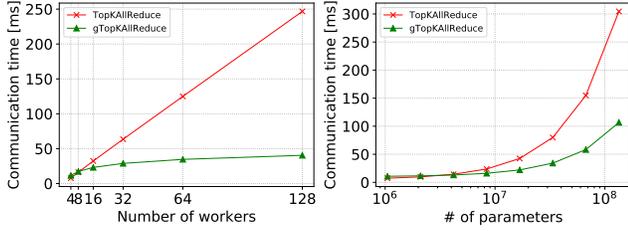[10]http://mvapich.cse.ohio-state.edu/benchmarks/

Fig. 9. Left: Time used for AllReduce algorithms on different number of workers with $m = 25 \times 10^6$ parameters, and the density of $\rho = 0.001$. Right: The time cost with respective to the number of parameters on 32 workers.

dense S-SGD has worst scaling efficiency because the full size of gradients makes the communication very slow on 1GbE clusters. The Top-$k$ S-SGD achieves some improvement on a small number of workers than S-SGD, but it has an obvious performance decrease when scaling to 32 GPUs. However, our proposed algorithm gTop-$k$ S-SGD achieves much more stable scaling efficiency even on clusters with a larger number of GPUs. For example, when scaling to 32 GPUs, our proposed gTop-$k$ S-SGD achieves $6.7\times$ faster than dense S-SGD on average, and $1.4\times$ improvement on average compared to Top-$k$. Particularly, gTop-$k$ S-SGD is up to $12\times$ and $1.7\times$ than S-SGD and Top-$k$ S-SGD respectively on the AlexNet model. The performance improvement will increase with the increase of number of workers. Summary of the training throughput on
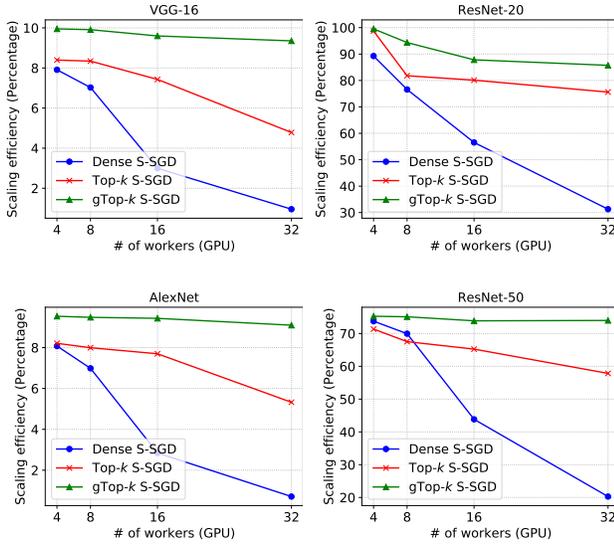


Fig. 10. Comparison of scaling efficiency of S-SGD with dense AllReduce (DenseAllReduce), Top-$k$ sparsification (TopKAllReduce) and gTop-$k$ sparsification (gTopKAllReduce), where $k = 0.001 \times m$. The higher the better.

different models is shown in Table. IV.

### E. Time performance analysis

We use the cases of 32 workers to analyze the time performance of gTop-$k$ S-SGD. We break down the time of an iteration into three parts: GPU computation time ($t_{compu.}$),

TABLE IV
THE SYSTEM TRAINING THROUGHPUT ON A 32-GPU CLUSTER.

| Model | Dense S-SGD | Top-$k$ | gTop-$k$ | $g/d$ | $g/t$ |
|---|---|---|---|---|---|
| VGG-16 | 403 | 2016 | 3020 | $7.5\times$ | $1.5\times$ |
| ResNet-20 | 9212 | 22272 | 25280 | $2.7\times$ | $1.1\times$ |
| AlexNet | 39 | 296 | 505 | $12.8\times$ | $1.7\times$ |
| ResNet-50 | 343 | 978 | 1251 | $3.65\times$ | $1.3\times$ |

Note: The throughput is measured with processed images per second. $g/d$ indicates the speedup of gTop-$k$ compared to the dense one, and $g/t$ indicates the speedup of gTop-$k$ compared to Top-$k$.

local sparsification time ($t_{compr.}$), and communication time ($t_{commu.}$). The results are shown in Fig. 11. On one hand, in the time breakdown of VGG-16 and AlexNet models, the communication overheads are much larger than computation because VGG-16 and AlexNet have three fully connected layers with a large number of parameters, while the computation is relatively fast. These also reflect that the scaling efficiency is low in Fig. 10 of S-SGD even with gTop-$k$ sparsification. On the other hand, the communication and sparsification overheads are much smaller than the computation with ResNet-20 and ResNet-50, which indicates low communication-to-computation ratios, so that the scaling efficiency can be up to $80\%$ even on the low-bandwidth network. Furthermore, it is noted that the time used by gradient sparsification is comparable to the computation time on VGG-16 and AlexNet models. The main reason is that Top-$k$ selection on GPU is inefficient and it could be non-trivial to be highly parallelized on SIMD architectures [32][33]. We will leave this as our future optimization direction.
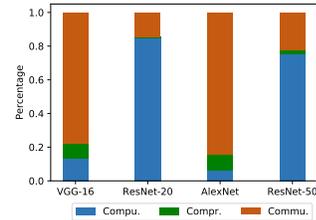


Fig. 11. Time breakdown of computation, compression and communication. "Compu." indicates forward and backward computation, "Compr." indicates the compression (sparsification) operation, and "Commu.' indicates gTop-$k$ gradients communication.

## V. DISCUSSION

### A. Convergence sensibility to the density

To understand the sensibility of the convergence to the density, we run the experiments with different values of the density $\rho$ using VGG-16 and ResNet-20 on the Cifar-10 data set on 4 workers. The convergence curves are shown in Fig. 12. It can be seen that even a very low density of $0.0005$ does not have a big impact on the model convergence to both models. However, a trade-off should be made to balance the high sparsification ratio and the convergence speed. One one hand, the higher sparsification would bring higher scaling efficiency to a larger number of workers. On the other hand,

one should also be careful to the upper bound of the sparsity that would hurt the model convergence.
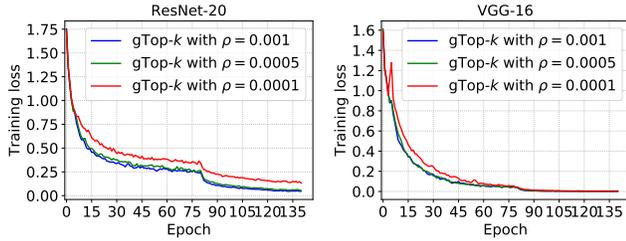


Fig. 12. Convergences with different $\rho$ on 4 workers.

*B. Convergence sensibility to the mini-batch size*

The previous section illustrates that gTop-$k$ S-SGD achieves nearly consistent convergences with proper choose densities. However, compared with Top-$k$ S-SGD, gTop-$k$ S-SGD has some disadvantages which may result in poorer generalization performance when the total number of iterations ($N$) is relatively small. Assume that we set $k = 0.001 \times m$, $P = 32$ and $B = b \times P = 1024$ on the Cifar-10 data set, we have $N = 5880$ with 120 epochs. At each iteration, Top-$k$ S-SGD could make $k \times P$ weights be updated, while gTop-$k$ S-SGD updates $k$ weights. Therefore, gTop-$k$ S-SGD has only 6 updates on some weights, while Top-$k$ S-SGD has about 192. The top-1 validation accuracy of VGG-16 and ResNet-20 trained with gTop-$k$ S-SGD and Top-$k$ S-SGD are shown in Fig. 13. It shows that ResNet-20 has ~9% accuracy degradation with gTop-$k$ S-SGD, while VGG-16 has only ~1% accuracy degradation with gTop-$k$ S-SGD. As a result, gTop-$k$ S-SGD requires more updates (by setting smaller mini-batch size) on ResNet-20 to achieve a higher accuracy, and it could also have higher accuracy degradation on VGG-16 by reducing the number of updates (by setting large mini-batch size). The comparison is shown in Fig. 14 with changed mini-batch sizes, which shows that gTop-$k$ S-SGD achieves closer accuracy (only ~0.5% degradation) to Top-$k$ S-SGD with smaller mini-batch size on ResNet-20, and gTop-$k$ has ~6% accuracy degradation with a large mini-batch size on VGG-16.
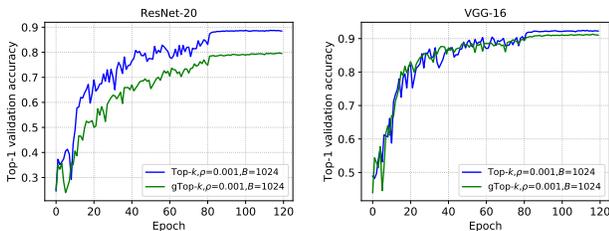


Fig. 13. The accuracy comparison between gTop-$k$ and Top-$k$ on ResNet-20 and VGG-16 with a mini-batch size of 1024 and $P = 32$.
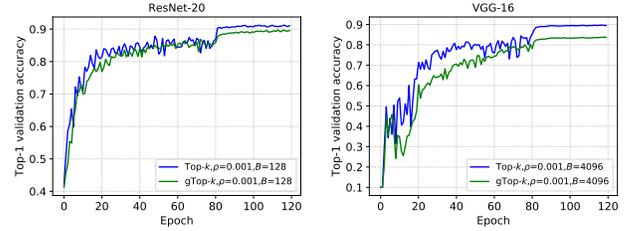


Fig. 14. The accuracy comparison between gTop-$k$ and Top-$k$ on ResNet-20 and VGG-16 with $P = 32$ and changed mini-batch sizes.

## VI. RELATED WORK

Gradient size reduction in communication is crucial for distributed synchronous SGD. Gradient quantization [34][35] and sparsification are two main techniques. Gradient quantization can only achieve a maximum of $32\times$ reduction compared to the 32-bit gradients, while gradient sparsification is more aggressive than quantization. Gradient sparsification zero-outs a large proportion of gradients to reduce the communication size dramatically. Aji et al. [17] and Chen et al. [11] empirically demonstrate that up to 99% gradients are not needed to update the model at each iteration, which indicates that the gradients would be very sparse to convergent the model with accumulations of gradient residuals. Aji et al. [17] use static threshold selection to determine $k$, while Chen et al. [11] propose a dynamic version. Lin et al. [12] further propose some optimization tricks (including the warmup strategy, momentum correction, and gradient clipping) to address the accuracy loss introduced by dropping a large number of gradients, and they show that Top-$k$ sparsification S-SGD can converge very close to S-SGD with dense gradients. The above techniques of quantization and sparsification can be combined to achieve a higher compression ratio of gradients with little accuracy loss. For example, Lin et al. [12] achieve up to $270\times$ and $600\times$ compression ratio without loss of accuracy. Researchers in [20] have realized that efficient sparse AllReduce algorithms are non-trivial to implement, and they propose the AllGather solution. However, the AllGather method requires a linear increase cost with respect to the number of workers. Therefore, the AllGather could be inefficient when scaling to large-scale clusters.

## VII. CONCLUSION AND FUTURE WORK

In this paper, we first showed that the accumulating results from top-$k$ gradients can be further sparsified by choosing largest absolute gradients before updating the model, which has no much impact on the model convergence. Then we identified that the Top-$k$ sparsification is inefficient in averaging the gradients from all workers because the indices of the Top-$k$ gradients are not the same such that one should use the AllGather collective to collect all the Top-$k$ gradients and indices. The AllGather method for Top-$k$ aggregation (TopKAllReduce) is linear expensive to the number of workers (i.e., the communication complexity is $O(kP)$, where

$P$ is the number of workers), so it would have very low scalability when scaling to large clusters. To this end, we proposed a global Top-$k$ (gTop-$k$) sparsification approach for S-SGD. The gradient aggregation algorithm based on gTop-$k$, named gTopKAllReduce, only requires a communication complexity of $O(k \log P)$. Experimental studies on various of deep neural networks including CNNs and RNNs were conducted to verify gTop-$k$ S-SGD has only slightly impact on the model convergence. The experiments conducted on the 32-GPU cluster inter-connected with 1 Gbps Ethernet showed that our proposed gTop-$k$ S-SGD has much higher scaling efficiency than S-SGD and Top-$k$ S-SGD.

Pipelining between computation and communication increases the scalability by optimally hiding the communication overheads [36]. In the future work, we would like to investigate layer-wise sparsification such that the communication overheads can be further overlapped by the computation tasks.

## VIII. ACKNOWLEDGEMENTS

## REFERENCES

[1] J. Dean, G. Corrado, R. Monga, K. Chen, M. Devin, M. Mao, A. Senior, P. Tucker, K. Yang, Q. V. Le *et al.*, "Large scale distributed deep networks," in *Advances in Neural Information Processing Systems*, 2012, pp. 1223–1231.

[2] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, p. 436, 2015.

[3] S. Shi, W. Qiang, and X. Chu, "Performance modeling and evaluation of distributed deep learning frameworks on GPUs," in *The Fourth International Conference on Big Data Intelligence and Computing*. IEEE, 2018.

[4] S. Shi, X. Chu, and B. Li, "A DAG model of synchronous stochastic gradient descent in distributed deep learning," in *The 24th International Conference on Parallel and Distributed Systems*. IEEE, 2018.

[5] D. Das, S. Avancha, D. Mudigere, K. Vaidynathan, S. Sridharan, D. Kalamkar, B. Kaul, and P. Dubey, "Distributed deep learning using synchronous stochastic gradient descent," *arXiv preprint arXiv:1602.06709*, 2016.

[6] P. Goyal, P. Dollár, R. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia, and K. He, "Accurate, large minibatch SGD: training ImageNet in 1 hour," *arXiv preprint arXiv:1706.02677*, 2017.

[7] W. Wang and N. Srebro, "Stochastic nonconvex optimization with large minibatches," *The 30th International Conference on Algorithmic Learning Theory*, 2019.

[8] Y. You, Z. Zhang, C.-J. Hsieh, J. Demmel, and K. Keutzer, "ImageNet training in minutes," in *The 47th International Conference on Parallel Processing*. ACM, 2018.

[9] X. Jia, S. Song, S. Shi, W. He, Y. Wang, H. Rong, F. Zhou, L. Xie, Z. Guo, Y. Yang, L. Yu, T. Chen, G. Hu, and X. Chu, "Highly scalable deep learning training system with mixed-precision: Training ImageNet in four minutes," *Workshop on Systems for ML and Open Source Software, collocated with NeurIPS 2018*, 2018.

[10] I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio, "Quantized neural networks: Training neural networks with low precision weights and activations," *The Journal of Machine Learning Research*, vol. 18, no. 1, pp. 6869–6898, 2017.

[11] C.-Y. Chen, J. Choi, D. Brand, A. Agrawal, W. Zhang, and K. Gopalakrishnan, "Adacomp: Adaptive residual gradient compression for data-parallel distributed training," in *The 32nd AAAI Conference on Artificial Intelligence*, 2018.

[12] Y. Lin, S. Han, H. Mao, Y. Wang, and W. J. Dally, "Deep gradient compression: Reducing the communication bandwidth for distributed training," in *International Conference on Learning Representations*, 2018.

[13] J. Wu, W. Huang, J. Huang, and T. Zhang, "Error compensated quantized SGD and its applications to large-scale distributed optimization," *International Conference on Machine Learning*, 2018.

[14] J. Bernstein, Y.-X. Wang, K. Azizzadenesheli, and A. Anandkumar, "SIGNSGD: Compressed optimisation for non-convex problems," in *International Conference on Machine Learning*, 2018, pp. 559–568.

[15] D. Alistarh, T. Hoefler, M. Johansson, N. Konstantinov, S. Khirirat, and C. Renggli, "The convergence of sparsified gradient methods," in *Advances in Neural Information Processing Systems*, 2018, pp. 5973–5983.

[16] S. U. Stich, J.-B. Cordonnier, and M. Jaggi, "Sparsified SGD with memory," in *Advances in Neural Information Processing Systems*, 2018, pp. 4452–4463.

[17] A. F. Aji and K. Heafield, "Sparse communication for distributed gradient descent," in *The 2017 Conference on Empirical Methods in Natural Language Processing*, 2017, pp. 440–445.

[18] P. Jiang and G. Agrawal, "A linear speedup analysis of distributed deep learning with sparse and quantized communication," in *Advances in Neural Information Processing Systems*, 2018, pp. 2530–2541.

[19] E. Chan, M. Heimlich, A. Purkayastha, and R. Van De Geijn, "Collective communication: theory, practice, and experience," *Concurrency and Computation: Practice and Experience*, vol. 19, no. 13, pp. 1749–1783, 2007.

[20] C. Renggli, D. Alistarh, and T. Hoefler, "SparCML: High-performance sparse communication for machine learning," *arXiv preprint arXiv:1802.08021*, 2018.

[21] J. Wangni, J. Wang, J. Liu, and T. Zhang, "Gradient sparsification for communication-efficient distributed optimization," in *Advances in Neural Information Processing Systems*, 2018, pp. 1299–1309.

[22] A. Sergeev and M. Del Balso, "Horovod: fast and easy distributed deep learning in TensorFlow," *arXiv preprint arXiv:1802.05799*, 2018.

[23] T. Hoefler, W. Gropp, R. Thakur, and J. L. Träff, "Toward performance models of MPI implementations for understanding application scaling issues," in *European MPI Users' Group Meeting*. Springer, 2010, pp. 21–30.

[24] S. Sarvotham, R. Riedi, and R. Baraniuk, "Connection-level analysis and modeling of network traffic," in *The 1st ACM SIGCOMM Workshop on Internet Measurement*. ACM, 2001, pp. 99–103.

[25] J. Pješivac-Grbović, T. Angskun, G. Bosilca, G. E. Fagg, E. Gabriel, and J. J. Dongarra, "Performance analysis of MPI collective operations," *Cluster Computing*, vol. 10, no. 2, pp. 127–143, 2007.

[26] A. Krizhevsky, V. Nair, and G. Hinton, "Cifar-10 (canadian institute for advanced research)," *URL http://www.cs.toronto.edu/kriz/cifar.html*, 2010.

[27] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *Computer Vision and Pattern Recognition*. IEEE, 2009, pp. 248–255.

[28] M. P. Marcus, M. A. Marcinkiewicz, and B. Santorini, "Building a large annotated corpus of English: The Penn Treebank," *Computational linguistics*, vol. 19, no. 2, pp. 313–330, 1993.

[29] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[30] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Computer Vision and Pattern Recognition*, 2016, pp. 770–778.

[31] A. Krizhevsky, "One weird trick for parallelizing convolutional neural networks," *arXiv preprint arXiv:1404.5997*, 2014.

[32] X. Mei, K. Zhao, C. Liu, and X. Chu, "Benchmarking the memory hierarchy of modern gpus," in *IFIP International Conference on Network and Parallel Computing*. Springer, 2014, pp. 144–156.

[33] A. Shanbhag, H. Pirk, and S. Madden, "Efficient Top-K query processing on massively parallel hardware," in *The 2018 International Conference on Management of Data*. ACM, 2018, pp. 1557–1570.

[34] 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 *Fifteenth Annual Conference of the International Speech Communication Association*, 2014.

[35] 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 *Advances in Neural Information Processing Systems*, 2017, pp. 1509–1519.

[36] S. Shi, X. Chu, and B. Li, "MG-WFBP: Efficient data communication for distributed synchronous SGD algorithms," in *INFOCOM 2019*, 2019.