
***ESync*: An Efficient Synchronous Parallel Algorithm for Distributed ML in Heterogeneous Clusters**

Zonghang LI

University of Electronic Science and
Technology of China
zhli0117@163.com

Huaman ZHOU

University of Electronic Science and
Technology of China
zhouhuaman@163.com

Hongfang YU

University of Electronic Science and
Technology of China
yuhf@uestc.edu.cn

Abstract

In recent years, researchers have paid more attention to distributed machine learning to train large-scale models with mass data in parallel, which is much more efficient than training on a single machine. Classical synchronous / asynchronous parallel algorithms (such as SSGD and ASGD) achieve excellent performance in aggregating model updates from devices in large-scale isomorphic clusters, however, these algorithms are not suitable in heterogeneous clusters, which consist of devices with different computing capabilities. This paper proposes *ESync*, an efficient synchronous parallel algorithm which allows devices in heterogeneous clusters to take full use of their computing capabilities while waiting for the straggling devices, meanwhile achieves the accuracy of SSGD and the training speed of ASGD with lowest traffic load. We evaluate *ESync* on a heterogeneous cluster with several Intel E5-2650v4 CPUs and GTX 1080TI GPUs, and find that *ESync* achieves $27\times$ training speedup to SSGD without loss of accuracy on AlexNet and achieves the state-of-the-art performance in both i.i.d. and non-i.i.d. data distributions.

1 Introduction

Massive data and large-scale machine learning models contributed a lot for the rapid development of Artificial Intelligence. In recent years, researchers have paid more attention to distributed machine learning to train large-scale models with mass data in parallel, which is much more efficient than training on a single machine. For example, Jia, Xianyan, et al.[1] trained ResNet-50 on ImageNet dataset with 2048 Tesla P40 GPUs and achieved 75.8% top-1 test accuracy in only 6.6 minutes. The training process of a standard distributed machine learning algorithm in synchronous mode can be divided into three steps: (1) Each device calculates model updates with local data independently; (2) Aggregate all the model updates and broadcast to all devices; (3) Update model parameters with aggregated model updates. Synchronous SGD (SSGD)[2] requires all the devices to aggregate model updates at the same pace, and devices with strong computing capabilities are blocked to wait for the straggling devices, which leads to waste of computational resources. Therefore, SSGD is preferred in isomorphic, stable and reliable clusters, in which devices have similar computing capabilities.

With the rapid development of GPU, TPU, FPGA, ASIC and other powerful computational resources, replacing all the old computing devices has become impractical and costly for enterprises and research institutes. Clusters are usually equipped with computing devices of different types and manufacturers, which leads to difference in computing capabilities among devices. Besides, Federated Learning[3]

supports cooperative training of millions of edge devices, such as mobile phones, smart watches, cameras and sensors, and the diversity of edge devices makes the heterogeneity of clusters much more significant.

The heterogeneity of clusters severely limits the training efficiency of SSGD to the weakest device, which causes the "Short Board" problem and results in inefficient training, low data throughput and great waste of computational resources. A natural idea to balance the computing time of each device in heterogeneous clusters is processing more samples while waiting for other devices, i.e. set larger batch size for devices with stronger computing capabilities. However, large batch size causes loss of accuracy[4]. In our preliminary experiments, we find that simply increase batch size according to the computing capability can not accelerate convergence, and probably lead to a decline in accuracy.

Asynchronous parallel algorithms (e.g. ASGD[5]) show strong adaptability to the heterogeneity of clusters (the performance reduction or failure of one device will not block other devices), but they also introduce "Delayed Gradient" problem, which uses outdated model updates to update the latest model, resulting in a decline in accuracy and an increase in the number of training iterations. Some researchers have analyzed the negative effect of delayed gradient to the convergence[6][7][8] and tried to improve ASGD from different perspectives to suppress the effect of delayed gradient on model accuracy, such as penalizing delay gradient by well-designed learning rate[9][10], compensate delay gradient[11], and control the difference of iterations[12]. However, these approaches will be difficult to meet the training and convergence requirements in clusters with strong heterogeneity, and ultimately lead to inefficient training and loss of accuracy.

To solve the "Short Board" problem and improve the training efficiency without loss of accuracy in heterogeneous clusters, this paper introduces *ESync*, an efficient synchronous parallel algorithm which allows devices in heterogeneous clusters to take full use of their computing capabilities while waiting for the stragglers. The key idea of *ESync* is to produce higher quality model updates that have a long-term vision while waiting for other devices. In each communication round, each device holds the same global model, saves an additional replica of the global model, and trains the replica for n_i iterations (on device i) via SGD (or other optimization algorithms) with local data. When all devices are ready, they calculate model updates based on difference between the global model and the latest replica, and send model updates to other devices to perform aggregation. Finally, the aggregated model updates are averaged and used to update the global model. The challenge of *ESync* is to assign the number of local iterations n_i to device i automatically according to the changing computing capabilities. To this end, *ESync* introduces State Server to coordinate the devices by ensuring that when the slowest device completes a local iteration, other devices have completed local iterations as many times as possible.

We evaluate *ESync* on a heterogeneous cluster with 2 Intel E5-2650v4 CPUs and 4 GTX 1080TI GPUs. Results show that *ESync* realizes a good balance of computing time among devices, achieves $7\times$ training speedup to SSGD on ResNet (up to $27\times$ on AlexNet) without loss of accuracy, processes $33\times$ more samples per second than SSGD, transmits $15\times$ less data per second than ASGD, and achieves the state-of-the-art performance in both i.i.d. and non-i.i.d. data distributions.

Contributions. To the best of the authors' knowledge, we propose the first efficient synchronous parallel algorithm to solve the "Short Board" problem in heterogeneous ML clusters. To sum up, the main contributions of this paper are as follows: (1) We introduce *ESync*, an efficient synchronous parallel algorithm that solves the "Short Board" problem, and greatly increases training speed without loss of accuracy; (2) We introduce State Server to assign the number of local iterations for each device automatically and dynamically, which provides fine-grained control over devices in dynamic clusters; (3) We provide a solution to take full use of the computing capabilities of the devices in heterogeneous clusters with lowest traffic load; (4) We provide a solution that allows aggregation operations to be performed in a synchronous manner in heterogeneous clusters, and supports many efficient collective communication algorithms (e.g. Ring Allreduce[13], Butterfly[14]).

The organization of this paper is as follows. In Section 2, we discuss the "Short Board" and "Delayed Gradient" problems and related work to improve them, and reveal the limitations of these work in heterogeneous clusters. In Section 3, we describe the *ESync* algorithm and analyze its convergence. In Section 4, we discuss the design of State Server and propose a design scheme. In Section 5, we evaluate *ESync* and compare its training speed, accuracy, data throughput, traffic load, computing time ratio with SSGD and ASGD. Finally, we conclude this paper and delineate several promising directions for future research in Section 6.

2 Related Work

This paper addresses the "Short Board" problem, and strives to take full use of the available computing capabilities of the devices in heterogeneous clusters, to accelerate the training of machine learning models significantly without loss of accuracy.

Short Board Problem. In Synchronous SGD (SSGD)[2], devices calculate model updates independently with local data, and push these model updates to other devices to perform aggregation. By using a barrier, devices with strong computing capabilities need to wait until other devices reach the barrier, which makes computational resources idle for a long time and causes the "Short Board" problem. To improve the training efficiency, Dean, Jeffrey, et al.[5] propose Asynchronous SGD (ASGD), which allows devices to perform computation and aggregation in different pace, with no need to wait for model updates from other devices. ASGD solves the "Short Board" problem and accelerates the training speed significantly. However, ASGD also introduces delayed gradient, which reduces the accuracy and convergence of the model, and the effect will be more serious when the difference in computing capabilities among devices is large in non-i.i.d. data distribution.

Delayed Gradient Problem. ASGD updates model $w_{t+\tau}$ with outdated model updates $g(w_t)$, resulting in a mismatch between the model updates and the model, and causes undesirable damage to the accuracy and the convergence of the model. Ho, Qirong, et al.[12] propose Stale Synchronous Parallel (SSP) to make a tradeoff between the accuracy and training speed by allowing the fastest device to be ahead of the slowest device by up to a bounded number of iterations τ . However, SSP still apply the outdated model updates to update the latest model and staleness still exist during the training process, and SSP will eventually degenerate into Bulk Synchronous Parallel (BSP)[15] with staleness in heterogeneous clusters when all the devices reach the threshold $t + \tau$ and wait for the slowest device, which makes the problems of "Short Board" and "Delayed Gradient" worse. McMahan[9] proposes AdaptiveRevision and Sra, Suvrit, et al.[10] propose Adadelay to penalize the delayed gradient by well-designed learning rate, but these approaches suppress the contributions of data on the stragging devices, which leads to unexpected loss of accuracy when the difference of steps is large. Zheng, Shuxin, et al.[11] propose DC-ASGD, which leverages the first-order term of Taylor expansion of the gradient function and approximates the Hessian matrix of the loss function to compensate delayed gradient, but its convergence requirements will be difficult to satisfy in clusters with strong heterogeneity, for the value of the high-order term $O((w_{t+\tau} - w_t)^2)I$ may be large and can not be ignored.

Instead of continuing to alleviate the negative effect of delayed gradient, *ESync* tackles the challenge of "Short Board" problem based on synchronous parallel algorithms directly to avoid the introduction of delayed gradient, and tries to utilize the blocking time of the aggregation operations to generate more valuable model updates.

3 The *ESync* Algorithm

Consider training a classification task on a heterogeneous cluster with p devices, and dataset are distributed on these devices $D = \{D_1, D_2, \dots, D_p\}$. For each dataset D_i on device i , we define the local loss l_i as:

$$l_i = \frac{1}{|D_i|} \sum_{(x_j, y_j) \in D_i} L(F(x_j, w), y_j) \quad (1)$$

where $|\cdot|$ denotes the size of the dataset, (x_j, y_j) denotes the data and label of the j -th sample in dataset D_i , $(F(\cdot), w)$ denotes the forward function and parameters of the model, and $L(\cdot)$ denotes the loss function. The global loss function on the distributed dataset can be defined as:

$$l = \frac{\sum_{i=1}^p |D_i| \cdot l_i}{|D|} \quad (2)$$

Our goal is to minimize the global loss l , i.e. to find the optimal shared parameters w^* :

$$w^* = \arg \min_w l_i \quad i = 1, 2, \dots, p \quad (3)$$

Gradient-descent techniques such as SGD is often used to solve such optimization problems. In this paper, we take mini-batch SGD as an example for simplicity, and it can be easily replaced by other optimization algorithms.

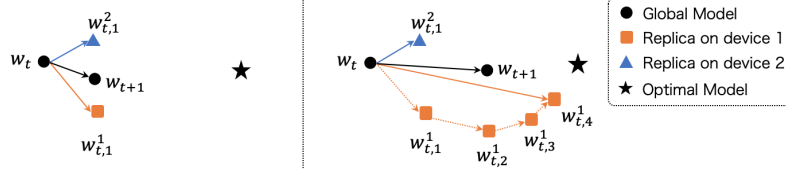


Figure 1: Overview of Sync (left) and *ESync* (right) within a communication round.

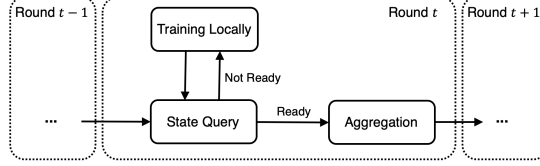


Figure 2: Training process of *ESync*.

3.1 Algorithm Modeling

Distributed machine learning tasks can be considered as a treasure hunt game with multiple explorers. In standard synchronous parallel algorithms (e.g. SSGD), each explorer shares one discovery in every step, and the explorers are forced to wait until all of them have a discovery. Instead, *ESync* allows experienced explorers to explore as many steps as possible, and share their discoveries when the freshmen are ready. An illustration of Sync and *ESync* within a communication round is shown in Figure 1, and the training process of *ESync* is shown in Figure 2.

We define the global model at the t -th communication round as w_t , and the replica of w_t at device i after k local iterations as $w_{t,k}^i$. w_0 is random initialized and shared by all devices. For $t > 0$, each device keeps the same global model w_t at the beginning, and saves a replica $w_{t,0}^i$ at device i :

$$w_{t,0}^i \leftarrow w_t \quad (4)$$

Step1: State Query. We introduce State Server to support state query, which decides whether devices should perform aggregation or continue training locally at the beginning of each local iteration. The State Server assigns the number of local iterations n_i for device i implicitly, where n_i is the number of times signal *NOT-READY* is received during a communication round.

Step2: Training Locally. For each local iteration k ($1 \leq k < n_i$), we sample a small batch data $D_e = \{(x_{i_1}, y_{i_1}), (x_{i_2}, y_{i_2}), \dots, (x_{i_b}, y_{i_b})\}$ from dataset D_i in random, calculate gradients and apply SGD to update the model replica $w_{t,k}^i$:

$$w_{t,k+1}^i \leftarrow w_{t,k}^i - \eta \nabla_w \frac{1}{b} \sum_{(x_j, y_j) \in D_e} L(F(x_j, w), y_j) \quad (5)$$

where η denotes the local learning rate, and b denotes the batch size. We can easily replace SGD with other optimization algorithms (e.g. Momentum, RMSProp, Adam) by modifying Formula (5).

Step3: Aggregation. When device i completes local iterations for n_i times, calculates the model updates based on difference between the latest replica w_{t,n_i}^i and the global model w_t :

$$\Delta w_t^i = w_{t,n_i}^i - w_t \quad (6)$$

then model updates Δw_t^i are sent to other devices to perform aggregation, and the aggregated model updates Δw_t are retrieved to each device to update the global model w_t :

$$\Delta w_t = \frac{1}{p} \sum_{i \in \{1, 2, \dots, p\}} \Delta w_t^i \quad (7)$$

$$w_{t+1} = w_t + \epsilon \Delta w_t \quad (8)$$

where ϵ denotes the global learning rate and can be simply set to 1.

Algorithm 1 *ESync* ($D_i, T, p, \eta, b, r, \epsilon = 1$)

Input: Dataset $D_i = \{(x_1, y_1), (x_2, y_2), \dots, (x_d, y_d)\}$; Maximum number of global iterations T ; Number of devices p ; Local learning rate η ; Batch size b ; Device id r ; Global learning rate ϵ (optional, the default value is 1).

- 1: *RANDOM-INITIALIZE-MODEL-PARAMETERS*(w_0, r);
- 2: *INITIALIZE-STATE-SERVER*(p, ϵ);
- 3: $c \leftarrow \text{RUN-BATCH-DATA-TEST}(b)$;
- 4: $t_e \leftarrow \text{GET-CURRENT-TIME}()$;
- 5: $t \leftarrow 0$; $k \leftarrow 0$;
- 6: $w_{0,0}^r \leftarrow w_0$;
- 7: **while** $t \leq T$:
- 8: $\text{ready} \leftarrow \text{QUERY-STATE}(r, k, c, t_e)$;
- 9: **if not** ready :
- 10: $t_s \leftarrow \text{GET-CURRENT-TIME}()$;
- 11: $\text{data} \leftarrow \text{GET-NEXT-RANDOM-MINIBATCH}(b)$;
- 12: $\text{gradients} \leftarrow \text{COMPUTE-AVERAGE-GRADIENTS}(w_{t,k}^r, \text{data})$;
- 13: $w_{t,k+1}^r \leftarrow w_{t,k}^r - \eta \cdot \text{gradients}$;
- 14: $t_e \leftarrow \text{GET-CURRENT-TIME}()$;
- 15: $c \leftarrow t_e - t_s$;
- 16: $k \leftarrow k + 1$;
- 17: **else**:
- 18: $\Delta w_t^r \leftarrow w_{t,k}^r - w_t$;
- 19: $\text{aggregated_updates} \leftarrow \text{GLOBAL-SUM-UPDATES}(r, \Delta w_t^r)$;
- 20: $w_{t+1} \leftarrow w_t + \epsilon \cdot \text{aggregated_updates} / p$;
- 21: $t \leftarrow t + 1$;
- 22: $w_{t,0}^r \leftarrow w_t$;
- 23: $k \leftarrow 0$;
- 24: $t_e \leftarrow \text{GET-CURRENT-TIME}()$;
- 25: $\text{RESET-STATE-SERVER}(r, t, t_e)$;

Output: Trained model parameters w_t .

The pseudo-code of *ESync* is shown in Algorithm 1. We can implement the aggregation operation *GLOBAL-SUM-UPDATES* at line 19 based on Parameter Server or MPI Allreduce, and apply compression technologies to reduce the amount of data that needs to be transmitted.

3.2 Convergence Analysis

In this section, we analyze the convergence of *ESync* and compare with the standard synchronous mini-batch SGD[16]. Assume the number of local iterations of p devices are $\{n_1, n_2, \dots, n_p\}$ ($n_i \geq 1, n_i \in \mathbb{Z}$). To simplify the analysis, we consider *ESync* as a synchronized algorithm that performs $\bar{n} = \frac{1}{p} \sum_{i=1}^p n_i$ local iterations on each device and $b\bar{n}p$ divides m without loss of generality.

Proposition 1. *Let $L(w, z)$ be a loss function with L -smooth convexity, and assume that the variance of the stochastic gradient $\nabla_w L(w, z)$ is bounded by σ^2 for all $w \in W$ and $z \in Z$. If the update rule has the serial expected optimality gap bound $\tilde{\psi}(\sigma^2, m) = \frac{1}{m} \psi(\sigma^2, m)$ when processed m samples, the expected optimality gap bound of Algorithm 1 is*

$$\tilde{\psi}\left(\frac{\sigma^2 \bar{n}}{bp}, \frac{m}{b\bar{n}p}\right) \quad (9)$$

specially, given $\psi(\sigma^2, m) = 2D^2L + 2D\sigma\sqrt{m}$, and Formula (11) is equal to

$$\frac{2b\bar{n}pD^2L}{m} + \frac{2D\sigma\bar{n}}{\sqrt{m}} \quad (10)$$

Proof of Proposition 2 is given in Appendix A. We can see that the bound of *ESync* is asymptotically equivalent to the bound of standard synchronous mini-batch SGD $O(\frac{1}{m} + \frac{1}{\sqrt{m}})$, which ensures that *ESync* can converge to the same accuracy as SSGD. Besides, *ESync* alleviates the "Large Batch Training" problem (larger batch size often leads to lower test accuracy)[4] by expanding the variance of the stochastic gradients by n_i times on each device.

4 State Server

We introduce State Server, which can be deployed on any devices, to assign the number of local iterations n_i for each device automatically and dynamically. The main idea is that, when the slowest device completes a local iteration, other devices have completed local iterations as many times as possible. The State Server is the key to solve the "Short Board" problem, and a good design for State Server can significantly reduce the blocking time of the aggregation operations.

4.1 State Query

Considering that computing capability may change, the State Server does not inform device i of its number of local iterations n_i directly. Instead, devices query the State Server whether they should perform aggregation or continue training locally in each iteration. This approach provides fine-grained control over devices in dynamic clusters, such as multi-tenant clusters and edge devices in Edge AI.

A simple but inefficient way is to perform aggregation when all devices complete at least one local iteration. Suppose there are two straggling devices A and B with the same computing capabilities. Although A and B complete computations at the same time, there will be a device that first reports to the State Server, and the device reports first (A) will continue training locally but the latter (B) starts to perform aggregation and then be blocked. Therefore, the "Short Board" problem is not solved because other devices still need to wait until A finishes its new iteration.

A better way is to perform aggregation when the querying device can not finish new iteration before the slowest device is ready. The pseudo-code of *QUERY-STATE* on State Server is shown in Algorithm 2. Assume the computing capabilities (time to process b samples) of p devices are $\{c_1, c_2, \dots, c_p\}$ ($c_1 \geq c_2 \geq \dots \geq c_p > 0$), and ε is a small constant to avoid boundary error. We can guarantee that the blocking time of the aggregation operation is less than c_p (0.03 seconds in our experiments), which is negligible compared to the slowest device (3.5 seconds in our experiments).

Algorithm 2 *State-Server*

```

1: procedure QUERY-STATE( $r, k, c, t_e$ ):
2:   // Record Format: ( $rank, local\_iterations, capability, last\_finish\_time, ready, round$ ).
3:   Update record current ( $rank\ r$ ) to ( $r, k, c, t_e, -1, -1$ ), where  $-1$  means unchanged;
4:    $slowest \leftarrow FIND-SLOWEST-DEVICE()$ ;
5:   if  $k == 0$  or  $slowest$  has not complete the PULL operation yet:
6:     return false;
7:    $current\_time \leftarrow GET-CURRENT-TIME()$ ;
8:    $rest\_time \leftarrow slowest.get\_capability() - current\_time + slowest.get\_last\_finish\_time()$ ;
9:   if  $current == slowest$  or  $slowest$  is ready or  $current.get\_capability() + \varepsilon > rest\_time$ :
10:     $current.set\_ready(true)$ ;
11:    return true;
12:  return false;

```

4.2 Non-Blocking Mode with Delayed ACK

Considering that the clusters may be in a network with high communication delay, and the blocked state query operation forces the devices to wait until the State Server responds signal *READY* or *NOT-READY*, which causes great waste of computational resources. Besides, the available network resources may be heterogeneous and dynamic, and the difference of communication delay introduces noise to the State Server.

We redefine capability c_i as the time to process b samples and the RTT of state query operation on device i to address both the difference of computing and communication capabilities, and allow devices to continue training locally while waiting for the response from the State Server to avoid waste of computational resources. We assume that implicit *NOT-READY* ACKs are received by default after sending state query requests, and devices continue training locally until an explicit *READY* ACK is received. When *READY* ACK of the k -th iteration is received, each device aborts the current $(k + \delta)$ -th iteration immediately and starts to perform aggregation based on iterations 1 to $k + \delta - 1$, where $\delta - 1$ is the number of local iterations that can be completed in RTT time.

5 Experiments

The experimental platform is built on MXNet[17], which is a flexible and efficient library for deep learning, and supports distributed deep learning by Parameter Server. We implement the aggregation operation on Parameter Server and an open source communication library PS-LITE. In this implementation, devices push model updates to the centralized KVStore to perform aggregation, and then pull aggregated model updates to local memory to update the model. We can also select MPI Allreduce algorithms to perform aggregation according to the number of devices, physical topology and size of the model. Besides, quantization and sparsification techniques can be used to accelerate communication, such as 1-bit SGD[18] and DGC[19]. We evaluate *ESync* on a heterogeneous cluster with 2 Intel E5-2650v4 CPUs and 4 GTX 1080TI GPUs, and compare its training speed, test accuracy, data throughput, traffic load, computing time ratio with standard synchronous parallel algorithm SSGD (Sync) and asynchronous parallel algorithm ASGD (Async).

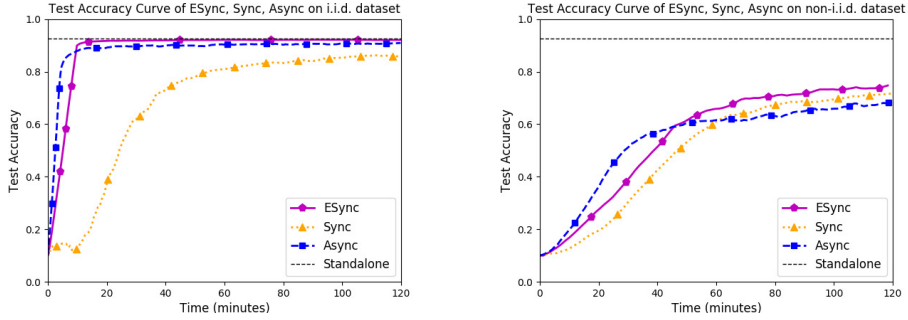


Figure 3: Test accuracy curve on i.i.d and non-i.i.d. Fashion-MNIST dataset.

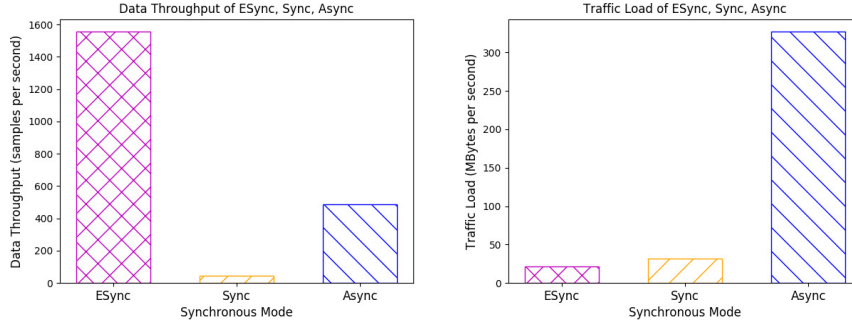


Figure 4: Data throughput and traffic load.

We set learning rate to 0.001 and batch size to 64, and apply mini-batch SGD for all experiments. We plot the curves of test accuracy over time, count the number of samples processed and the amount of data transmitted per second, and count the total time of computing and communication. Furthermore, we define the speedup as the time ratio to achieve test accuracy 0.8 that is close to convergent accuracy, and the computing time ratio as the proportion of total computing time to total training time.

Test accuracy and training speed. *ESync* greatly increase training speed without loss of accuracy. We distribute the classic dataset Fashion-MNIST[20] to each device in i.i.d. and non-i.i.d. data distributions and train ResNet[21] respectively. Results show that *ESync* achieves the state-of-the-art performance in both settings. As shown in Figure 3, while training on i.i.d. dataset, *ESync* achieves almost the same training speed as Async, which is about $7\times$ faster than Sync, and achieves the test accuracy of standalone training in the shortest time. Furthermore, while training on non-i.i.d. dataset, *ESync* converges faster than Sync without loss of accuracy while the accuracy of Async is impaired.

Data throughput and traffic load. *ESync* enables heterogeneous clusters to process data more efficiently, meanwhile minimizes traffic load in network. As shown in Figure 4, *ESync* trains $33\times$ more samples than Sync and $3\times$ more samples than Async per second, however, only $1/15$ of the data need to be transmitted per second compared to Async. The reason is that *ESync* performs multiple iterations locally, which increases data throughput and reduces the number of communication.

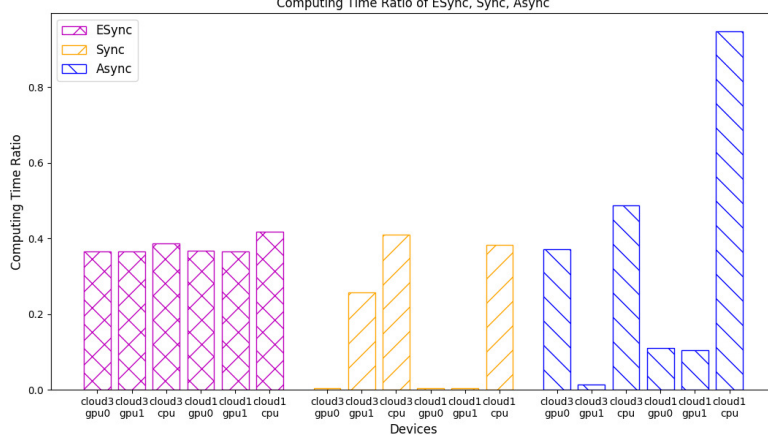


Figure 5: Computing time ratio on each device.

Computing time. *ESync* improves the utilization rate of computing power on devices with strong computing capabilities significantly. As shown in Figure 5, Sync spends too much time waiting for straggling devices and Async spends little time on computing but always in communication. *ESync* realizes a good balance of computing time among devices, enables devices with strong computing capabilities to make full use of their computational resources while waiting for the straggling devices.

We compare the speedup of *ESync* to Sync and the test accuracy with the standalone training on several classic models, as shown in Table 1. The results show that *ESync* can greatly accelerate the training process and finally achieve the same (sometimes higher) accuracy as the standalone training.

Table 1: Speedup and accuracy of *ESync* on different models.

Model	Speedup	Test Accuracy	
		<i>ESync</i>	Standalone
AlexNet	27×	0.929	0.928
Inception-v3	7×	0.937	0.932
ResNet-18-v1	7×	0.926	0.926
ResNet-50-v1	4×	0.914	0.910
ResNet-50-v2	6×	0.918	0.920
MobileNet-v1	2×	0.903	0.902

6 Conclusion

This paper proposes *ESync*, an efficient synchronous parallel algorithm which allows devices in heterogeneous clusters to take full use of their computing capabilities while waiting for the straggling devices, meanwhile achieves the accuracy of SSGD and the training speed of ASGD with lowest traffic load. We introduce State Server to assign the number of local iterations for each device automatically and dynamically, aim to solve the "Short Board" problem in heterogeneous and dynamic clusters. Specially, *ESync* allows aggregation operations to be performed in a synchronous manner in heterogeneous clusters, and provides users with flexibility in selecting different collective algorithms and communication topologies according to the characteristics of tasks and network (e.g. star topology of Parameter Server, ring topology of MPI Ring Allreduce, hierarchical hybrid topologies, and even user-defined topologies).

In future work, we are going to design a communication mechanism based on *ESync* that supports user-defined communication mode (including communication topologies and algorithms), and propose an algorithm for automatic communication mode design based on the characteristics of clusters (e.g. number of devices, computing resources, communication resources), and realize seamless switching of communication modes while detecting dramatic changes in the clusters to better adapt to heterogeneous and dynamic clusters, such as multi-tenant clusters and edge devices in Edge AI.

References

- [1] Jia, Xianyan, et al. "Highly scalable deep learning training system with mixed-precision: Training imagenet in four minutes." arXiv preprint arXiv:1807.11205, 2018.
- [2] Chen, Jianmin, et al. "Revisiting distributed synchronous SGD." arXiv preprint arXiv:1604.00981 (2016).
- [3] Konecny, Jakub, et al. "Federated learning: Strategies for improving communication efficiency." arXiv preprint arXiv:1610.05492 (2016).
- [4] Keskar, Nitish Shirish, et al. "On large-batch training for deep learning: Generalization gap and sharp minima." arXiv preprint arXiv:1609.04836 (2016).
- [5] Dean, Jeffrey, et al. "Large scale distributed deep networks." Advances in neural information processing systems. 2012.
- [6] Lian, Xiangru, et al. "Asynchronous parallel stochastic gradient for nonconvex optimization." Advances in Neural Information Processing Systems. 2015.
- [7] Avron, Haim, Alex Druinsky, and Anshul Gupta. "Revisiting asynchronous linear solvers: Provable convergence rate through randomization." Journal of the ACM (JACM) 62.6 (2015): 51.
- [8] Agarwal, Alekh, and John C. Duchi. "Distributed delayed stochastic optimization." Advances in Neural Information Processing Systems. 2011.
- [9] McMahan, Brendan, and Matthew Streeter. "Delay-tolerant algorithms for asynchronous distributed online learning." Advances in Neural Information Processing Systems. 2014.
- [10] Sra, Suvrit, et al. "Adadelat: Delay adaptive distributed stochastic convex optimization." arXiv preprint arXiv:1508.05003 (2015).
- [11] Zheng, Shuxin, et al. "Asynchronous stochastic gradient descent with delay compensation." Proceedings of the 34th International Conference on Machine Learning-Volume 70. JMLR. org, 2017.
- [12] Ho, Qirong, et al. "More effective distributed ml via a stale synchronous parallel parameter server." Advances in neural information processing systems. 2013.
- [13] Gibiansky, Andrew. "Bringing HPC techniques to deep learning". <http://research.baidu.com/bringing-hpc-techniques-deep-learning/>, 2017.2.21.
- [14] Rabenseifner, Rolf. "Optimization of collective reduction operations." International Conference on Computational Science, Krakow, Poland, 2004.6.6-6.9.
- [15] Gerbessiotis, Alexandros V., and Leslie G. Valiant. "Direct bulk-synchronous parallel algorithms." Journal of parallel and distributed computing 22.2 (1994): 251-267.
- [16] Dekel, Ofer, et al. "Optimal distributed online prediction using mini-batches." Journal of Machine Learning Research 13.Jan (2012): 165-202.
- [17] Chen, Tianqi, et al. "Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems." arXiv preprint arXiv:1512.01274 (2015).
- [18] Seide, Frank, et al. "1-bit stochastic gradient descent and its application to data-parallel distributed training of speech dnns." Fifteenth Annual Conference of the International Speech Communication Association. 2014.
- [19] Lin, Yujun, et al. "Deep gradient compression: Reducing the communication bandwidth for distributed training." arXiv preprint arXiv:1712.01887 (2017).
- [20] Xiao, Han, Kashif Rasul, and Roland Vollgraf. "Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms." arXiv preprint arXiv:1708.07747 (2017).
- [21] He, Kaiming, et al. "Deep residual learning for image recognition." Proceedings of the IEEE conference on computer vision and pattern recognition. 2016.
- [22] Krizhevsky, Alex, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." Advances in neural information processing systems. 2012.
- [23] Szegedy, Christian, et al. "Rethinking the inception architecture for computer vision." Proceedings of the IEEE conference on computer vision and pattern recognition. 2016.
- [24] Howard, Andrew G., et al. "Mobilenets: Efficient convolutional neural networks for mobile vision applications." arXiv preprint arXiv:1704.04861 (2017).
- [25] Xiao, Lin. "Dual averaging methods for regularized stochastic learning and online optimization." Journal of Machine Learning Research 11.Oct (2010): 2543-2596.

A Proof of Proposition 1

Proposition 2. Let $f(w, z)$ be a loss function with L -smooth convexity, and assume that the variance of the stochastic gradient $\nabla_w L(w, z)$ is bounded by σ^2 for all $w \in W$ and $z \in Z$. If the update rule has the serial expected optimality gap bound $\tilde{\psi}(\sigma^2, m) = \frac{1}{m}\psi(\sigma^2, m)$ when processed m samples, the expected optimality gap bound of Algorithm 1 is

$$\tilde{\psi}\left(\frac{\sigma^2 \bar{n}}{bp}, \frac{m}{b\bar{n}p}\right) \quad (11)$$

specially, given $\psi(\sigma^2, m) = 2D^2L + 2D\sigma\sqrt{m}$, and Formula (11) is equal to

$$\frac{2b\bar{n}pD^2L}{m} + \frac{2D\sigma\bar{n}}{\sqrt{m}} \quad (12)$$

Proof. Let \bar{D}_j denote the set of batches used in all local iterations at j -th communication round on all devices. In each communication round, we define z_{s+lb}^k as the $(s + lb)$ -th sample and w_l^k as the model at l -th local iteration on device k , and assume that each batch is sampled i.i.d. from D_k on device k . Let $\bar{f} : W \times \bar{D}_j \rightarrow \mathbb{R}$ be defined as

$$\bar{f}(w, (z_1^1, z_2^1, \dots, z_{bn_p}^p)) = \frac{1}{bp} \sum_{k=1}^p \sum_{l=1}^{n_k} \sum_{s=1}^b f(w_l^k, z_{s+lb}^k) \quad (13)$$

using the linearity of the gradient operator, we have

$$\nabla_w \bar{f}(w, (z_1^1, z_2^1, \dots, z_{bn_p}^p)) = \frac{1}{bp} \sum_{k=1}^p \sum_{l=1}^{n_k} \sum_{s=1}^b \nabla_w f(w_l^k, z_{s+lb}^k) \quad (14)$$

where p denotes the number of devices, b denotes the batch size for each iteration, n_k denotes the number of local iterations on device k at current communication round, and $(z_1^1, z_2^1, \dots, z_{bn_p}^p)$ denotes the sequence of samples in \bar{D}_j .

Let $\bar{z}_j = (z_1^1, z_2^1, \dots, z_{bn_p}^p)$ and $F(w) = \mathbb{E}_{\bar{z} \in \bar{D}_j} \bar{f}(w, \bar{z}) = \mathbb{E}_{z \in D} f(w, z)$, and \bar{f} is L -smooth and convex as well as f due to the triangle inequality. We write the variance of the stochastic gradient using Euclidean norm as

$$\| \nabla_w \bar{f}(w, \bar{z}_j) - \nabla F(w) \|^2 \quad (15)$$

$$= \left\| \frac{1}{bp} \sum_{k=1}^p \sum_{l=1}^{n_k} \sum_{s=1}^b (\nabla_w f(w_l^k, z_{s+lb}^k) - \nabla F(w_l^k)) \right\|^2 \quad (16)$$

$$= \frac{1}{b^2 p^2} \sum_{k=1}^p \sum_{l=1}^{n_k} \sum_{s=1}^b \sum_{s'=1}^b \left\langle \nabla_w f(w_l^k, z_{s+lb}^k) - \nabla F(w_l^k), \nabla_w f(w_l^k, z_{s'+lb}^k) - \nabla F(w_l^k) \right\rangle \quad (17)$$

we have z_s and $z_{s'}$ ($s \neq s'$) are independent since they are sampled i.i.d., therefore

$$\mathbb{E} \left\langle \nabla_w f(w_l^k, z_{s+lb}^k) - \nabla F(w_l^k), \nabla_w f(w_l^k, z_{s'+lb}^k) - \nabla F(w_l^k) \right\rangle \quad (18)$$

$$= \left\langle \mathbb{E}[\nabla_w f(w_l^k, z_{s+lb}^k) - \nabla F(w_l^k)], \mathbb{E}[\nabla_w f(w_l^k, z_{s'+lb}^k) - \nabla F(w_l^k)] \right\rangle = 0 \quad (19)$$

then, let $\bar{n} = \frac{1}{p} \sum_{i=1}^p n_i$, we have the variance bound for every $w \in W$,

$$\| \nabla_w \bar{f}(w, \bar{z}_j) - \nabla F(w) \|^2 \quad (20)$$

$$= \frac{1}{b^2 p^2} \sum_{k=1}^p \sum_{l=1}^{n_k} \sum_{s=1}^b \mathbb{E} \| \nabla_w f(w_l^k, z_{s+lb}^k) - \nabla F(w_l^k) \|^2 \quad (21)$$

$$\leq \frac{1}{b^2 p^2} \sum_{k=1}^p \sum_{l=1}^{n_k} \sum_{s=1}^b \sigma^2 = \frac{\sigma^2 \bar{n}}{bp} \quad (22)$$

formula (22) shows that the variance of the stochastic gradient $\nabla_w \bar{f}(w, \bar{z}_j)$ is bounded by $\frac{\sigma^2 \bar{n}}{bp}$ for every $w \in W$.

Here we define the initial model at j -th communication round as \bar{w}_j^0 (shared by all devices). If the update rule has a regret bound $\psi(\sigma^2, m)$ over m inputs, the regret bound for \bar{f} over $m/b\bar{n}p$ communication rounds is

$$\mathbb{E}[R(m)] = \mathbb{E} \left[\sum_{j=1}^{m/b\bar{n}p} (\bar{f}(\bar{w}_j^0, \bar{z}_j) - \bar{f}(w^*, \bar{z}_j)) \right] \leq \psi\left(\frac{\sigma^2 \bar{n}}{bp}, \frac{m}{b\bar{n}p}\right) \quad (23)$$

according to Theorem 3 by Xiao[25], for convex loss functions and i.i.d. inputs, the optimality gap $G(m)$ always have

$$\mathbb{E}[G(m)] \leq \frac{1}{m} \mathbb{E}[R(m)] \quad (24)$$

therefore, for the specific regret bound $\psi(\sigma^2, m) = 2D^2L + 2D\sigma\sqrt{m}$, we have

$$\mathbb{E}[G(m)] \leq \tilde{\psi}\left(\frac{\sigma^2\bar{n}}{bp}, \frac{m}{b\bar{n}p}\right) = \frac{1}{m} \psi\left(\frac{\sigma^2\bar{n}}{bp}, \frac{m}{b\bar{n}p}\right) \leq \frac{2b\bar{n}pD^2L}{m} + \frac{2D\sigma\bar{n}}{\sqrt{m}} \quad (25)$$

□