

GossipGraD: Scalable Deep Learning using Gossip Communication based Asynchronous Gradient Descent

Jeff Daily

Pacific Northwest National Laboratory
jeff.daily@pnnl.gov

Abhinav Vishnu

Pacific Northwest National Laboratory
abhinav.vishnu@gmail.com

Charles Siegel

Pacific Northwest National Laboratory
charles.m.siegel@gmail.com

Thomas Warfel

Pacific Northwest National Laboratory
thomas.warfel@pnnl.gov

Vinay Amatya

Pacific Northwest National Laboratory
vinay.amatya@pnnl.gov

ABSTRACT

In this paper, we present *GossipGraD* – a gossip communication protocol based Stochastic Gradient Descent (SGD) algorithm for scaling Deep Learning (DL) algorithms on large-scale systems. The salient features of GossipGraD are: 1) reduction in overall communication complexity from $\Theta(\log(p))$ for p compute nodes in well-studied SGD to $O(1)$, 2) model diffusion such that compute nodes exchange their updates (gradients) indirectly after every $\log(p)$ steps, 3) rotation of communication partners for facilitating direct diffusion of gradients, 4) asynchronous distributed shuffle of samples during the feedforward phase in SGD to prevent over-fitting, 5) asynchronous communication of gradients for further reducing the communication cost of SGD and GossipGraD. We implement GossipGraD for GPU and CPU clusters and use NVIDIA GPUs (Pascal P100) connected with InfiniBand, and Intel Knights Landing (KNL) connected with Aries network. We evaluate GossipGraD using well-studied dataset ImageNet-1K (≈ 250 GB), and widely studied neural network topologies such as GoogLeNet and ResNet50 (current winner of ImageNet Large Scale Visualization Research Challenge (ILSVRC)). Our performance evaluation using both KNL and Pascal GPUs indicates that GossipGraD can achieve perfect efficiency for these datasets and their associated neural network topologies. Specifically, for ResNet50, GossipGraD is able to achieve $\approx 100\%$ compute efficiency using 128 NVIDIA Pascal P100 GPUs – while matching the top-1 classification accuracy published in literature.

1 INTRODUCTION

Deep Learning (DL) algorithms are a class of Machine Learning and Data Mining (MLDM) algorithms. A deep neural network (DNN) – which stores the model of a DL algorithm – contains several *layers of neurons* inter-connected with *synapses*. By using a large number of layers, DL algorithms are able to conduct transformations on highly non-linear data which is a common property in many science datasets. DL algorithms have demonstrated remarkable improvements in many Computer Vision tasks [1, 2] and science domains such as High Energy Physics [3], and Climate Modeling [4]. Several DL implementations such as TensorFlow [5], Caffe [6], Theano [7, 8], and Torch [9] have become available.

At the same time, DL algorithms are undergoing a tremendous revolution of their own. Several types of DL algorithms – primarily geared by input properties (tabular, images, speech) – are used by researchers/practitioners. As an example, Multi-layer Perceptrons (MLPs) are widely used for tabular data. Similarly for images,

videos and speech, Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) are the *de facto* choice. However, CNNs and RNNs are very computationally expensive, requiring days of training time using GPUs even on modest size datasets. Their computational requirements are further worsened by: 1) very deep neural networks such as GoogLeNet [10] and Residual Networks (*ResNet*) [11](Figure 1), and 2) an increasing volume of data produced by simulations, experiments and handheld devices.

An important solution to these problems is distributed memory DL algorithms that are capable of execution on multi-device (such as multi-GPUs) and large scale systems. Table 1 shows the prominent approaches for distributed memory DL implementations. We observe that the HPC ready implementations typically incur $\Theta(\log(p))$ communication complexity, while a few of the HPC ready implementations support overlap of communication with computation.

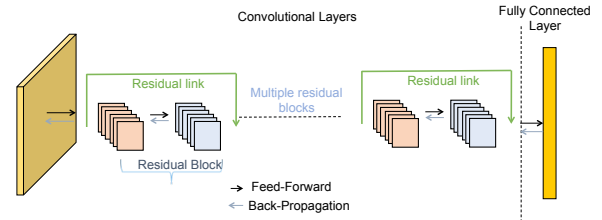


Figure 1: An example of a neural network – ResNet [11]. ResNet has two types of layers – convolutional and fully connected. Besides feedforward and backprop links, it has a residual link. The residual block is a unit which is repeated to create very deep topologies

DL algorithms primarily use the *Gradient Descent* method, which is an iterative technique to update the weights of *synapses* using the *error* between the *ground truth* (actual value) and the *predicted value* (using the current state of the DNN). The larger the difference, the steeper GD converges to a minima – a low value of minima generates the solution. An important type of Gradient Descent is Batch/Stochastic Gradient Descent (SGD) – where a random subset of samples are used for iterative *feed-forward* (calculation of predicted value) and *back-propagation* (update of synaptic weights).

Consider a strong scaling scenario where a batch (b) of the input dataset is split across multiple compute nodes (p) – an example of *data parallelism*. Under such a scenario, the computation is bounded by $\Theta(\frac{b}{p})$, although the communication – requiring an all-to-all reduction – is bounded by $\Theta(\log(p))$. With weak scaling – the

work per compute node remains constant, but the communication increases by $\log(p)$. This argument is further validated by Goyal *et al.* [12] and PowerAI [13], since weak scaling beyond it results in accuracy loss. Additionally, in practice the actual communication time deviates from $\Theta(\log(p))$ due to system issues as pointed out by other researchers [14, 15].

Name	Comm. Complexity	HPC Ready	Comm. Overlap	Leverage Bisection BW
FireCaffe [16]	$\Theta(\log(p))$	✓	✗	✓
S-Caffe [17]	$\Theta(\log(p))$	✓	✓	✓
MaTeX [18]	$\Theta(\log(p))$	✓	✗	✓
Poseidon [19]	$O(1)$	✗	✓	✗
Petuum [20]	$O(1)$	✗	✗	✗
GeePS [21]	$O(1)$	✗	✓	✗
ProjectAdam [22]	$O(1)$	✗	✓	✗
TensorFlow [5]	$O(1)$	✗	✗	✗
MXNET [23]	$O(1)$	✓	✗	✗
CaffeonSpark [24]	$O(1)$	✗	✓	✗
SparkNet [25]	$O(1)$	✗	✓	✗
CNTK [26]	$\Theta(\log(p))$	✓	✗	✓
Parle [27]	$O(1)$	✗	✗	✗
PaddlePaddle [28]	$O(1)$	✓	✓	✗
DistBelief [29]	$O(1)$	✓	✗	✗
Caffe2 [12]	$\Theta(\log(p))$	✓	✓	✓
DSSTNE [30]	$O(1)$	✓	✗	✗
Chainer [31]	$\Theta(\log(p))$	✓	✗	✓
PowerAI [13]	$\Theta(\log(p))$	✓	✓	✓
EASGD3 [32]	$O(\log(p))$	✓	✓	✗
Blot <i>et al.</i> [33]	$O(1)$	✗	✓	✗
Others [34–36]	$O(1)$	✗	✓	✗
GossipGraD	$O(1)$	✓	✓	✓

Table 1: Comparison of GossipGraD with the major distributed Deep Learning approaches. HPC ready implementations leverage the HPC interconnects natively such as either using MPI or native interfaces.

Several researchers have proposed methods to alleviate the communication requirements of distributed SGD [21, 23, 37–39]. The parameter server based approaches (shown in Figure 2(a)) use a server to hold the latest copy of the model, while worker(s) send gradients (computed as a function of error) and request the latest weights. An important aspect of the parameter server is the expected constant communication complexity ($O(1)$ instead of $\Theta(\log(p))$) for SGD). However, researchers have pointed out the deficiencies of parameter server based approaches [17, 39]: 1) a single parameter server becomes a performance bottleneck, 2) multiple parameter servers (such as using GPUs) result in waste of compute resources [17], 3) communication between parameter servers in presence of multiple servers is expensive, and 4) parameter servers require warm-start to facilitate convergence [39]. Because of those limitations, parameter servers have only been practical for relatively small clusters, and are therefore excluded from further consideration in this paper.

Data parallelism based implementations such as FireCaffe [16], S-Caffe [17], PowerAI [13], Caffe2 [12], MaTeX [18], and others [40] use all-to-all reduction and other collective operations and use the bisection bandwidth of communication network effectively (table 1).

Recently proposed S-Caffe [17], PowerAI [13] and Caffe2 [12] also overlap communication with compute (back-propagation). However, even with these optimizations, communication becomes a bottleneck at scale even with asynchronicity such as evident in their evaluation.

Hence, alternate – possibly complimentary – approaches are required to scale DL implementations further. There are two primary approaches published in literature of gossip based DL as proposed by Jin *et al.* [41] and Blot *et al.* [33]. Under these approaches, a random communication partner is selected for sending model updates as shown in Figure 2(b). These updates are applied by the receiving compute node and the process is repeated iteratively. However, both Jin *et al.* and Blot *et al.* report performance and convergence degradation at modest scale. Up on careful review, we attribute the problems to: 1) communication imbalance since random selection does not guarantee balanced communication across compute nodes, 2) communication overhead due to lack of asynchronicity, and 3) poor convergence at scale due to imbalanced diffusion of gradients. Architecture specific approaches such as proposed by You *et al.* [32] provide support for GPU and KNL architectures, but their efficiency decreases sharply on using 64 KNL nodes – which implies that it is not a feasible solution for larger scale systems.

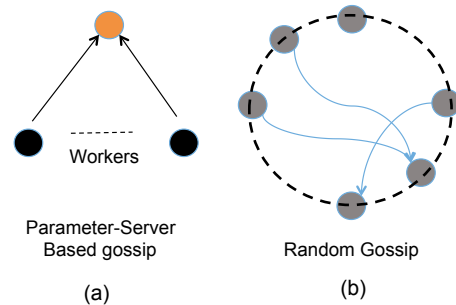


Figure 2: (a) Parameter-server based implementation – which is an extreme form of all-to-one gossip, (b) Random gossip as considered by Blot *et al.* [33] and Jin *et al.* [41]

1.1 Problem Definition and Synopsis

Existing all-to-all reduction based approaches provide excellent convergence, but limited scalability. At the same time, parameter server and existing gossip communication based approaches provide constant communication complexity, but limited convergence. The important question is: *Can we design a set of distributed memory DL algorithms that provides constant communication complexity while leveraging asynchronous communication, but similar convergence properties as sequential SGD?* That is the objective of this paper.

To achieve this objective, we propose GossipGraD – a novel gossip communication protocol based SGD – to design distributed DL algorithms. Our algorithm addresses the limitation of existing gossip based approaches by eliminating communication imbalance, introducing asynchronicity such that the communication is completely overlapped with computation, and shuffling the communication partners and the training dataset asynchronously to provide better diffusion of model updates and preventing over-fitting.

1.2 Contributions

Specifically, we make the following contributions in this paper:

- We present the design choices for communication partner selection in GossipGraD. Specifically, we present partner selection such that each compute node communicates with exactly one partner at each step and after $\log(p)$ steps, all compute nodes have communicated indirectly.
- We present techniques for batch-wise rotation of partners for further enabling direct diffusion of model updates. We consider layer-wise/batch-wise approaches for asynchronous communication of model updates to minimize the observed communication time during the training phase.
- To prevent over-fitting (a situation where the model has memorized the training set), we present *asynchronous* distributed memory shuffle of samples. The shuffle does not incur additional communication complexity since it is overlapped with the feedforward step.
- We provide theoretical underpinning of GossipGraD (and its heuristics), and present the case that GossipGraD converges to a similar solution as default SGD.
- We use Caffe – a widely available toolbox – as our implementation baseline. We implement GossipGraD using vendor optimized Caffe – Intel-Caffe for Knights Landing (KNL) based systems and NVIDIA-Caffe for GPU based systems. Each of these implementations are also extended for distributed memory implementation by using MPI [42, 43] such as considered by other researchers including Caffe2, PowerAI, and S-Caffe.

We evaluate GossipGraD using the well-studied dataset ImageNet-1K (1.3M images, $\approx 250\text{GB}$) [44], and neural network topologies such as GoogLeNet [2] and ResNet [11] using 128 (32 x 4) Pascal P100 GPUs, and LeNet3 and CIFARNet using 32 node Intel KNL system (Aries interconnect). Our performance evaluation indicates the effectiveness of GossipGraD on a variety of architectures. Specifically, we are able to provide complete overlap of communication with computation for ResNet50 on 128 GPUs resulting in $> 99\%$ efficiency and match top-1 accuracy published by other researchers (PowerAI, Caffe2 and Chainer). As an example, after 30 just training epochs of ResNet50, GossipGraD achieves a top-1 accuracy of 50% – which matches with the accuracy published by Caffe2, Chainer and PowerAI. Our theoretical underpinnings and accuracy evaluation indicate that GossipGraD provides similar accuracy as the well-studied SGD algorithm. Hence, GossipGraD is suitable for deployment on extreme scale systems.

The rest of the paper is organized as follows: In section 2, we present the background of the proposed research. We present baseline setup in section 3, GossipGraD design in section 4, introduce asynchronicity in section 5 and theoretical proof of convergence in section 6. We discuss experimental results in section 7, related work in section 8, followed by conclusions in section 9.

2 BACKGROUND

2.1 Gradient Descent Algorithm

The gradient descent algorithm – which iteratively updates the weights of a deep neural network to a local minimum of the cost

function – is the most widely used DL algorithm. In the algorithm, each sample (which could be tabular observation, an image, or speech/text), is an input to the *feed-forward* step. The output of each feed-forward step is a *predicted value*, which is compared with the *label* (ground truth), and their difference is used to calculate the *gradients* – which are applied for updating the weights. In essence, the feed-forward step calculates the *loss* – a measure of difference between label and predicted values, and the objective is to minimize the loss on the training set, while ensuring that the validation loss decreases as well. The general technique is typically referred to as *gradient descent* towards a minimum of the loss function (Figure 3).

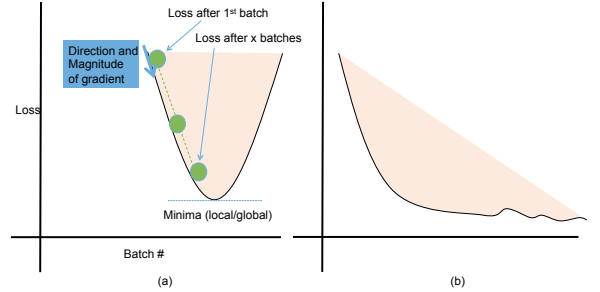


Figure 3: (a) A pictorial representation of a loss curve as function of batch – each batch implies a model update, (b) A typical DL loss curve containing long-tail and non-convex shape

For gradient descent and its variants, the rule(s) for updating the weights during the back-propagation step can be presented as follows:

$$\mathbf{w}' = \mathbf{w} + \lambda \nabla_{\mathbf{w}} C \quad (1)$$

$$\mathbf{b}' = \mathbf{b} + \lambda \nabla_{\mathbf{b}} C \quad (2)$$

where \mathbf{w} are the weights, \mathbf{b} the biases, λ the learning rate, and C is a cost function to be optimized which is usually square error or cross-entropy. To compute the derivatives, we set $W^{(\ell)}$, $b^{(\ell)}$ to be the weights and biases for each layer, $z^{(\ell+1)} = W^{(\ell)} a^{(\ell)} + b^{(\ell)}$ and $a^{(\ell)} = \sigma(z^{(\ell)})$, where σ is the activation function, and we set n_{ℓ} represent the number of layers. The details of the algorithm are presented in Algorithm 1.

Algorithm 1 Back Propagation [45]

- 1: **input:** Data $X \in \mathbb{R}^{n \times p}$ and labels $Y \in \mathbb{R}^{n \times l}$
 - 2: Compute all $z^{(\ell)}$ and $a^{(\ell)}$.
 - 3: $\delta^{(n_{\ell})} = -(y - a^{(n_{\ell})}) \odot \sigma'(z^{(n_{\ell})})$
 - 4: **for** ℓ from $n_{\ell} - 1$ to 1 **do**
 - 5: $\delta^{(\ell)} = W^{\ell} \delta^{(\ell+1)} \odot \sigma'(z^{(\ell)})$
 - 6: **end for**
 - 7: $\nabla_{W^{(\ell)}} C = \delta^{(\ell+1)} a^{(\ell)T}$
 - 8: $\nabla_{b^{(\ell)}} C = \delta^{(\ell+1)}$
-

3 BASELINE SETUP

In this section, we present a solution space for scaling SGD on distributed memory systems. We first present a distributed memory

	Method	Symbol
1	Input Dataset	X
2	Batch size	b
3	Number of Processes	p
4	Network Latency	l
5	Network Bandwidth	$\frac{1}{G}$
6	Set of Layers	$L = L_0..L_{n-1}$

Table 2: Symbols used for GossipGraD and Other Approaches

implementation of *Vanilla SGD* (simply referred as SGD for rest of the paper). SGD provides a baseline for performance and accuracy comparisons with GossipGraD. Table 2 shows the symbols used for time-space complexity analysis of the proposed approaches.

3.1 Distributed Memory SGD

These are several design choices for scaling SGD on distributed memory systems. Specifically, to maintain equivalence to the sequential algorithm, a strong scaling approach may be useful. Considering a batch size (b), the overall work for (p) processes is expected to be $\Theta\left(\frac{b}{p} + \|L\| \log p\right)$, where $\|L\|$ is the overall size of the synaptic weights in a DNN, and an all-to-all reduction is conducted among all MPI processes. The $\log(p)$ term is obtained due to the time-complexity of the reduction tree typically implemented as a binomial/k-nomial tree in distributed systems. For weak-scaling, we expect the time complexity to be $\Theta(b + \|L\| \log p)$, since the work/process remains constant. Figure 4 shows an example of distributed memory SGD with ResNet topology.

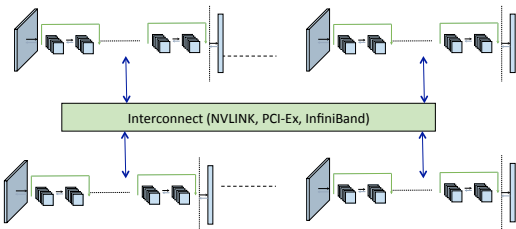


Figure 4: Pictorial representation for SGD using ResNet example. Data parallelism is used by replicating the ResNet model across all compute nodes. All-to-all reduction is executed after back-propagation is complete.

3.2 Scaling Limitations of SGD: Possible Solutions

An advantage of distributed SGD is its strict equivalence to the sequential SGD algorithm. However, the communication complexity – bounded by $\Theta(\log(p))$ – becomes a limiting factor at scale. Hence, it is important to consider alternate techniques to scale SGD on large scale systems. One possibility is to communicate to fewer processes – possibly to exactly one process after each batch update. This reduces the communication complexity from $\Theta(\log(p))$ to $O(1)$. In section 4, we consider design approaches for GossipGraD which achieves this objective.

Another (complimentary) possibility is to introduce asynchronicity in SGD (AGD) and GossipGraD. Recently proposed solutions such as S-Caffe [17], Caffe2 [12] and PowerAI [13] overlap communication with computation during the back-propagation step. This is evident from Algorithm 1, since gradients for each layer are ready for communication before back-propagation is initiated on preceding layers. We present design choices for asynchronous execution in section 5.

4 GOSSIPGRAD DESIGN

In this section, we present a solution space for designing GossipGraD. We first present an intuition for GossipGraD. We consider several elements of GossipGraD: 1) communication partner selection for exchanging model updates, 2) techniques for asynchronous data shuffle for preventing over-fitting, 3) low-overhead partner rotation for better diffusion of model updates and accelerating convergence to the solution, and 4) implementations on state of the art GPU and CPU clusters.

4.1 Extreme Case: No Communication

Let us consider the basic SGD algorithm presented earlier. After each batch, an all-to-all reduction (typically implemented using `MPI_Allreduce`) is conducted, which takes $\log(p)$ communication steps for completion. The objective of GossipGraD is to reduce the communication complexity of SGD from $\log(p)$ to $O(1)$. One possibility is to completely eliminate communication from SGD. This approach is attractive for minimizing communication complexity. However, the no-communication approach has two problems: 1) each model (by data parallelism) learns only on a subset of data. This implies that the models drift apart further at increasing scale of compute nodes, and 2) this approach would produce an ensemble of DL models, while our objective is to produce a single model at the end of the training phase. Hence, we disregard no communication as a viable solution to addressing the communication complexity of DL algorithms.

4.2 Intuition for Gossip Communication

Gossip – as frequently observed in social networks – is frequently used in computer-to-computer communications especially in distributed systems. Compute nodes may select a random partner for gossip, and each exchange is not expected to be reliable. However – over a period of time – gossip communication is expected to provide robust information across partners. This property of gossip communication protocols is the premise of GossipGraD. As each pair of compute nodes gossip their model updates – with a potentially different partner at each step – the expectation is that individual models of compute nodes steadily converge towards the same model (theoretical underpinnings are presented in section 6). SGD usually takes several epochs (as an example ResNet requires 90 epochs), where each epoch may have thousands of batches. The premise of GossipGraD is that by *gossiping* between compute nodes over thousands of batches, the model updates are diffused across all compute nodes.

Figure 2 demonstrated the limitations of existing gossip based approaches, including the parameter server approaches, which can be considered as an extreme form of gossip. The approaches presented

by Jin *et al.* [41] and Blot *et al.* [33] which use random gossip suffer from communication imbalance, poor diffusion of gradients, lack of data shuffle and asynchronous exchange of gradients. The proposed GossipGraD design intends to address each of these limitations. We begin with a discussion on selection gossip communication partner.

4.3 Gossip Communication Partner Selection

The approach of selecting a communication partner should have the following properties: 1) constant communication complexity – such that the solution is scalable to very large scale systems, 2) balanced communication such that at each step there is a unique set of communicating nodes, 3) diffusion of gradients in sub-linear time, and 4) leveraging the bisection bandwidth of the communication network.

To facilitate constant communication complexity and balanced communication, we propose a hierarchical virtual organization of compute nodes such as in hypercube topology (shown in Figure 5(a)). Other configurations such as dissemination-based (shown in Figure 5(b)) have similar property as hypercube virtual topology. However, dissemination-based approach is better since it communicates with exactly two compute nodes at each step. Another important property of hypercube based and dissemination based approaches is that all compute nodes have communicated indirectly with each other in exactly $\log(p)$ batches – which achieves our objective of sub-linear diffusion of model updates.

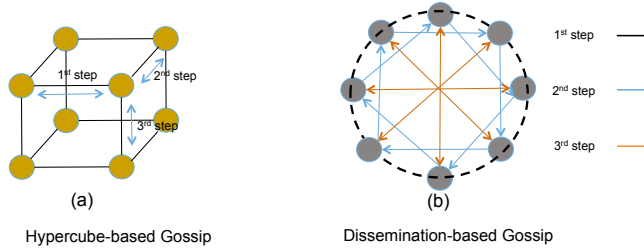


Figure 5: (a) Hypercube based virtual topology organization of compute nodes (b) dissemination based organization of compute nodes

4.4 Expected diffusion of gradients with Hypercube and Dissemination Virtual Topologies

4.4.1 Diffusion in Hypercube Topology. An example is shown in Figure 6. Intuitively, the expectation is that by exchanging gradients with a partner, the gradients are diffused *indirectly* over a period of steps. As evident from the figure, after $\log(p)$ steps, the gradients have been indirectly diffused with all compute nodes.

4.4.2 Diffusion in Dissemination Topology. The primary difference between dissemination and hypercube algorithms is the partner selection. At step k of the all-to-all reduction, a process p_i sends data to a process with MPI rank $(p_i + 2^k)\%p$, and receives from $(p_i + p - 2^k)\%p$. Figure 7 shows the steps in classical all-to-all reduction and GossipGraD, when using the dissemination algorithm. Similar to hypercube exchange, dissemination provides $\log(p)$ time complexity.

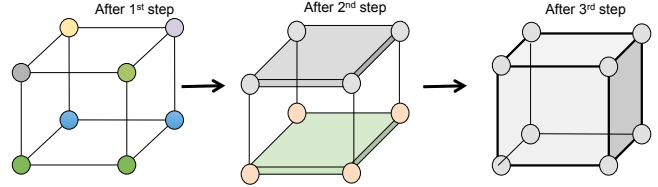


Figure 6: The expected diffusion of gradients across steps. After 1st step, set of pairs have same gradient, after 2nd step, top and bottom surfaces, and after 3rd step, the entire set of processes.

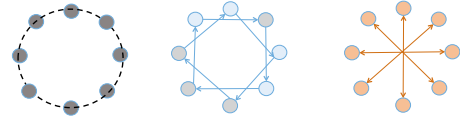


Figure 7: An example of exchanges in dissemination algorithm for GossipGraD. Unlike hypercube exchange, each process sends and receives from a different partner and completes indirect exchange in $\log(p)$ steps

However, for GossipGraD, the difference between dissemination and hypercube based partner selection is worth consideration. Specifically, by using the dissemination algorithm, GossipGraD is diffusing gradients from two partners at each step, while hypercube exchange is diffusing from exactly one partner. Hence, we primarily consider dissemination exchange based partner selection for GossipGraD.

4.5 Handling Side-Effects of Reduced Communication

A possible side-effect of the reduced communication is divergence from the sequential algorithm – since only a small subset of compute nodes are communicating at each step. We propose two techniques to address this:

4.5.1 Partner Rotation. Let us consider an execution of dissemination exchange based GossipGraD. We observe that the communication partners are repeated after $\log(p)$ steps. As a result, the direct diffusion of gradients is restricted to a small fraction $(\frac{\log(p)}{p})$ partners. To alleviate this limitation, we propose a *partner rotation* based approach. The premise of the rotation based approach is such that the communication partners are modified after every $\log(p)$ steps. The rotation based approach facilitates direct diffusion of gradients to all compute nodes over a period of time, without increasing the time complexity of GossipGraD. To achieve this objective, we consider p random shuffles of the original communicator. After each $\log(p)$ steps, the next shuffled communicator is used for creating a new virtual dissemination topology. Since the communicators are created at start of the application, the overall cost of creating them is easily amortized over long running training phase of the DL implementation.

4.5.2 Rotation of Samples. The majority of distributed memory DL implementations read the samples once in memory at the

beginning of the training phase, and use them repeatedly till convergence. This approach is sufficient for default SGD implementation – since all compute nodes communicate after every batch. Since GossipGraD reduces the communication complexity, we propose a distributed *shuffle* of the samples such that after a compute node has used a batch of samples, it sends the samples to another compute node. To develop a different topology from the dissemination topology of gradients, we consider a ring virtual topology for the shuffle of samples, where each compute node sends a recently completed batch to its neighbor. A simple yet effective shuffle provides a nice property that each sample is considered again for feedforward/backpropagation step by a compute node only after all other compute nodes have considered it once.

5 ASYNCHRONOUS GOSSIPGRAD

As discussed in the previous section, we have reduced the communication complexity from $\Theta(\log(p))$ to $O(1)$. However, we have added the shuffle of samples – potentially increasing the overall communication observed by the compute nodes at the training phase. Another important consideration is leveraging asynchronous communication for exchanging model updates. We leverage the property of feedforward and backpropagation step to asynchronously shuffle the samples and exchange the model updates.

An important observation from Algorithm 1 and GossipGraD design is that the gradients for each layer are ready for diffusion before the gradients for preceding layers are computed. Hence, it is possible to overlap the computation of gradients for preceding layers by communicating gradients of the current layer. Existing approaches such as S-Caffe [17], PowerAI [13] and Caffe2 [12] have made similar observation for all-to-all reduction. We use a similar property for point-to-point communication in GossipGraD.

5.1 Distributed Asynchronous GossipGraD using MPI Non-Blocking Primitives

For GossipGraD, we use non-blocking point-to-point communication primitives – `MPI_Isend` and `MPI_Irecv`. We implement GossipGraD by initiating non-blocking send/receive, as soon as the gradients for a layer are ready. Figure 8 shows a pictorial representation. The non-blocking MPI requests return an MPI handle. To maximize the available overlap, our implementation generates the communication requests and executes `MPI_TestAll` followed by an `MPI_WaitAll` function, after the gradients for all the layers have been computed.

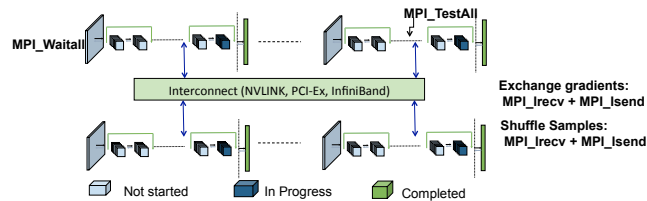


Figure 8: A representative execution of AGD and GossipGraD using non-blocking MPI primitives

5.2 Distributed Asynchronous GossipGraD using Asynchronous Progress Thread

A critical problem with non-blocking point-to-point and collective operations is that it requires MPI runtimes to make asynchronous progress/hardware support. However, this feature is not necessarily available in all implementations. At the same time, blocking operations (such as `MPI_Allreduce`) are heavily optimized by vendors, since they are used in many large scale applications. Let us consider ResNet50 – which has about 25M parameters (100 MBytes of data) – to motivate the need for an asynchronous thread based implementation. Many of the layers communicate large data ($> 1\text{MBytes}$) – which is much greater than the size of *rendezvous threshold* in MPI implementations. Hence, we expect that the majority of communication uses rendezvous protocol for point-to-point communication, which necessitates asynchronous communication progress. Sur *et al.* have proposed MPI implementations to support asynchronous progress [46] – although they may not be available in practice.

Hence, we use an asynchronous thread to ensure progress on blocking point-to-point and collective communication primitives. Since each message transfer is sufficiently large to saturate the network bandwidth, a single communication thread is sufficient. Once the gradients for a layer are generated, the main compute thread *enqueues* a pointer to the gradients on a *communication queue*. The asynchronous thread dequeues from the communication queue and issues a blocking point-to-point/all-to-all reduction operation. Figure 9 provides a pictorial representation of the asynchronous implementation.

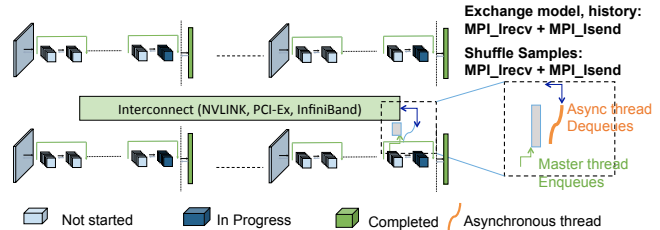


Figure 9: GossipGraD implementation using blocking collectives with asynchronous communication threads. Main thread enqueues the gradients, which are dequeued by the communication thread

5.2.1 Considerations for Implementation. There are several important considerations in this implementation: 1) The thread contention is insignificant, since the master thread enqueues the gradients, and a single asynchronous thread, which dequeues for point-to-point/all-to-all reduction, 2) the implementation requires multi-threaded MPI which is available on most platforms, 3) presence of an asynchronous thread does not affect the performance of GPU based implementation (since thread is scheduled on CPUs), for CPU based implementation a single communication thread negligibly affects the available parallelism.

Yet, in practice the performance of multi-threaded MPI is not consistent across all platforms. Specifically, we have observed a performance degradation when `MPI_THREAD_MULTIPLE` is used, and depending on the implementation, an asynchronous progress

	Meaning	Symbol
1	Weights and Biases	w
2	Cost Function	C
3	i th sample	x_i
4	i th label	y_i^L
5	predicted value for i th sample	y_i

Table 3: Symbols used for Proof

may not be provided. Fortunately, using `MPI_TestAll` – which is a non-blocking operation to invoke progress engine in MPI – we have observed an expected overlap of communication with computation. Hence, we discard the asynchronous thread based implementation and use the non-blocking primitive and `testall` based implementation for communication.

6 SKETCH OF PROOF OF CONVERGENCE

In this section, we provide a sketch of a proof that the GossipGraD algorithm converges to a similar model as SGD at the end of training phase. We use the symbols described in Table 3. Recall that for a classification network, the *cost function* is the cross-entropy function defined by

$$C(w) = -\frac{1}{n} \sum_{i=1}^n y_i^L \cdot \log(y_i)$$

Let the number of nodes be p and assume that the data is distributed between nodes.

LEMMA 6.1. *With shuffling, the compute nodes have the same cost function.*

PROOF. Without shuffling, the cost functions are summations over the samples resident on the node. However, shuffling ensures that over time, each sample is only re-considered for feedforward and backpropagation once the compute node has considered every other sample. Hence, the cost function being optimized is the summation over all samples, and thus identical across nodes. \square

Due to the structure of neural networks, the cost function $C(w)$ is non-convex (this non-convexity is analyzed in detail in [47]). The non-convexity implies that there are several possible local minima even when executed on a single device. Even though the cost functions are same across all compute nodes, they may end up at different local minima due to the stochastic nature of SGD. Hence, communication of model updates is imperative to guarantee convergence to a single local minimum. Naturally, this problem becomes worse at scale.

In section 3, we have presented communication approaches where after each batch, the model on each compute node is averaged with another unique compute node. Using the communication approaches presented earlier, we obtain:

THEOREM 6.2. *Each compute node in GossipGraD converges to a local minimum of the cost function.*

PROOF. The key technical result is [48] which shows that for any process (called a *positive supermartingale*) where the expectation value of the $n + 1^{th}$ batch, conditioned on the first n batches, is nonnegative and at most the value of the n th batch, converges to a

limit with probability 1. This result has been extended in [49, 50] to show that SGD converges (with probability 1) for any cost function that is continuous and which is differentiable at all but finitely many points.

By Lemma 6.1, each compute node is optimizing a copy of the same objective function to find a local minimum. We follow the standard proof for SGD, where our goal is to show that the sequence $\{C^k = C(w_k)\}_{k=0}^{\infty}$ is a positive supermartingale. Positivity follows from the structure of $C(w)$, an average of explicitly positive terms, such as the negative of the logarithms of probabilities. We use the properties of SGD described below to show that the expectation value of the cost is at most the value of the cost in a given batch.

The properties of SGD that are typically used to show that it is a supermartingale are: 1) the weight differences are bounded, 2) the set of non-minima where the gradient vanishes is encountered with probability zero, and 3) the expectation of the difference in weights before and after each batch points towards a local minimum (denoted by w^*), that is,

$$\mathbb{E}[(w_{n+1} - w_n) \cdot (w_{n+1} - w^*)] \geq 0.$$

The first two properties hold for GossipGraD because they hold for SGD, so only the third property remains. Here, the difference between GossipGraD and SGD is that w_{n+1} for SGD is defined by global gradient update, but $w_{n+1,j}$ for GossipGraD on compute node j is defined by $(W_{n+1,j} + W_{n+1,c_i(j)})/2$ where $W_{i,j}$ are the weights at batch i on compute node j after gradient updates and c_i is the permutation at batch i determining which other compute node is being averaged into compute node j .

As n increases, the probability that this expectation value is negative decreases, such that for any $\epsilon > 0$, there exists an N such that for all $n > N$, we have

$$\Pr(\mathbb{E}[(w_{n+1} - w_n) \cdot (w_{n+1} - w^*)] < 0) < \epsilon.$$

The desired inequality does not hold only in pathological cases, where a compute node j has consistently moved away from the local minimum. As in the case of SGD, however, the expectation of the gradients points towards the local minimum. Hence, the desired inequality holds with probability 1 as the number of batches increase, such as observed in practical cases of large datasets and DNNs. \square

COROLLARY 6.3. *All compute nodes converge to the same local minimum of the cost function.*

PROOF. Theorem 6.2 states that the model on each compute node converges to some local minimum. We prove that they converge to the same local minima by using contradiction. Let us assume that all models have converged, but not to the same local minimum. Then, for some batch, models that are at different local minimum are averaged together due to partner rotation. However, this would cause the weights of these models to change, contradicting the assumption that they have already converged. Therefore, all models must have converged to the same local minimum. \square

7 PERFORMANCE EVALUATION

In this section, we present an in-depth evaluation of GossipGraD and its associated heuristics using datasets such as MNIST, CIFAR10 and well-studied neural network topologies on ImageNet-1K dataset. The table 4 provides a description of the architectures used for

performance evaluation. Table 5 provides a brief description of the datasets used for performance evaluation. Table 6 provides a description of the heuristics.

7.1 Baseline Setup

DL algorithms have a few hyperparameters which are setup at the start of the training phase. These include the batch size, learning rate and momentum. We fix the momentum to be exactly same as provided by default versions of Intel-Caffe and NVIDIA-Caffe. Since we use weak-scaling, the effective batch size is multiplied by the number of devices (such as number of GPUs or number of compute nodes). For handling weak-scaling with SGD/AGD, we use suggestions by Krizhevsky [56] that multiplies the learning rate on a single device by \sqrt{p} . For GossipGraD, we use the same batch size on each compute node as defined for a single device (such as a single GPU or a single KNL compute node), and keep the learning rate unchanged. We use AGD – a theoretically equivalent asynchronous gradient descent implementation – as suggested by PowerAI, Caffe2 and Chainer as the baseline for performance and accuracy comparisons. We consider AGD to be a better baseline than default SGD since AGD provides better performance than the SGD, and provides theoretical equivalence to SGD.

7.2 Comparing Gossip to AGD: Evaluation on (LeNet3) MNIST and (CIFARNet) CIFAR10

7.2.1 Performance Results. Figures 10 and 11 show the relative speedup of GossipGraD on MNIST and CIFAR10 datasets to AGD implementation on KNL and P100 clusters respectively. MNIST on 32 GPUs needs 1.2s per epoch (batch size on each device is 64), while CIFAR10 on 32 GPUs needs 0.75s per epoch (batch size on each device is 100). Both MNIST and CIFAR10 are relatively small datasets, but provide a delicate balance between communication (since network sizes are small) and computation (since compute is relatively small). We observe the following: 1) relative speedup of GossipGraD to AGD for P100 on both MNIST and CIFAR10 is higher in comparison to KNL, since a single P100 GPU is much faster than single KNL node, and 2) even with weak-scaling, the relative speedup continues to increase since GossipGraD is able to overlap the communication effectively due to $O(1)$ communication, while AGD is unable to do so even in the presence of a layer-wise asynchronous communication.

7.2.2 Accuracy Results. Figure 12 shows the accuracy charts for GossipGraD on KNL, GossipGraD on P100 and AGD implementation as a function of number of epochs. We use 32 P100 GPUs (8 compute nodes) and 32 KNL compute nodes for the evaluation to validate whether the accuracies match at the largest scale of evaluation on MNIST considered in this paper. We observe that the validation accuracy saturates to $\approx 99.2\%$ for three implementations – which empirically proves the argument that GossipGraD provides similar accuracy as AGD implementation. Figure 13 compares the accuracy charts for GossipGraD on KNL, GossipGraD on P100 and AGD implementation using 32 KNL nodes and 32 GPUs on CIFAR10 dataset. The three implementations roughly track each other and converge towards similar accuracy ($\approx 72-73\%$), which is within the margin of error.

7.2.3 Lessons Learned. We observe that Gossip provides similar accuracy while providing significant speedup in comparison to the AGD implementation even on relatively small scale KNL and P100 clusters. We observe about 1.9x speedup for MNIST since GossipGraD is able to overlap the communication of both samples and model updates with the feedforward and backpropagation step. We also conclude that GossipGraD is able to completely overlap the communication – which implies that it is useful for small size datasets as well.

7.3 Comparing GossipGraD to AGD: ResNet 50 Evaluation

7.3.1 Performance Comparisons. Table 7 shows the compute efficiency of GossipGraD and compares it with published efficiency on PowerAI [13] using up to 128 P100 GPUs. We use a batch size of 32 on each device for GossipGraD implementations – which is the batch size used by other researchers. The GossipGraD implementation use asynchronous `MPI_TestAll` based implementation. Since we use weak-scaling, the overall batch size on 128 GPUs is 4096. We observe that the overall compute efficiency of GossipGraD is $\approx 100\%$ independent of the number of GPUs. The PowerAI implementation achieves 100% efficiency for 4 and 8 GPUs but continues to decrease to 95% on 128 GPUs. At 32 batch size per device for GossipGraD, the feedforward and backpropagation time is about 96ms (providing 10.4 batch updates/second). The synchronous point-to-point communication time is 27ms which is completely overlapped by the GossipGraD implementation.

7.3.2 Accuracy Comparisons. Figure 14 shows the validation accuracy as a function of epochs for ResNet50 using GossipGraD on 128 P100 GPUs. ResNet50 uses a step learning training regimen, which is executed for 100 epochs. After every 30 epochs, the current learning rate is multiplied by 0.1. Since ResNet50 is well studied by other researchers including PowerAI [13], Caffe2 [12] and Chainer [31], we compare the GossipGraD validation accuracy after a fixed number of epochs with these approaches published elsewhere.

For GossipGraD, we use the original learning rate provided for ResNet50 (0.1) and follow the training regimen presented earlier. We compute the validation accuracies after every ≈ 5 epochs. At about 30 epochs, we see an accuracy of 50%. This which is the accuracy reported by Chainer, PowerAI and Caffe2 for 30 epochs. Hence, GossipGraD provides similar accuracy as well-published literature while providing nearly perfect overlap of communication with computation.

7.4 Comparing GossipGraD with AGD: Evaluating GoogLeNet

Figure 15 shows the relative speedup of GossipGraD in comparison to AGD using up to 32 GPUs for GoogLeNet respectively. We use a batch size of 16 on each compute node for both GossipGraD and AGD. We observe: 1) relative speedup for GossipGraD increases with scale – even on weak scaling – since the overall communication time increases. GoogLeNet has 5M parameters (20M floats) – which is much smaller than ResNet which has 25M parameters. However, ResNet50 is relatively computationally expensive than GoogLeNet.

Name	CPU (#cores)	# GPUs/node	Baseline	Network	MPI	Nodes	#cores	# GPUs
P100	Power8 (20)	NVIDIA Pascal P100 (4)	NVIDIA-Caffe[51]	NVLink, IB-EDR	IBM MPI	32	N/A	128
KNL	Intel KNL (68)	N/A	Intel-Caffe[52]	Cray Aries	Cray MPICH	32	2176	N/A

Table 4: A description of system architecture and associated software

Dataset	Neural Network	Description	Training Samples	Validation Samples	Image Size	Classes
MNIST [53]	LeNet3 [54]	Handwritten Digits	60000	10000	28×28	10
CIFAR-10 [55]	CIFARNet	Small Images	60000	10000	$32 \times 32 \times 3$	10
ImageNet	GoogLeNet [2]	Diverse Images	1281167	50000	$256 \times 256 \times 3$	1000
ImageNet	ResNet [11]	Diverse Images	1281167	50000	$256 \times 256 \times 3$	1000

Table 5: A description of datasets and associated neural network topologies

Name	Type	Implemented	Description of DL Algorithm and Implementation
AGD	Asynchronous GD	Yes	Implements SGD by asynchronous layer-wise communication.
GossipGraD	Batch-wise Gossip GD with Rotation of partners and samples	Yes	Implements batch-wise AGD using Gossip and Rotation of Samples

Table 6: A description of GossipGraD Approaches. We implement these approaches and evaluate them on GPU and CPU systems.

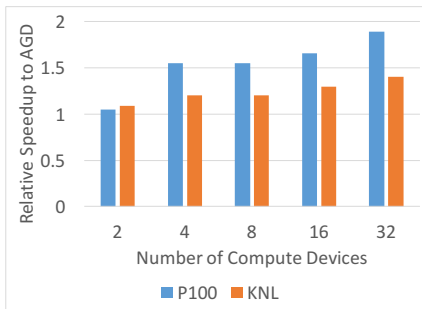


Figure 10: Relative Speedup of Gossip to AGD for P100 and KNL Clusters on MNIST dataset

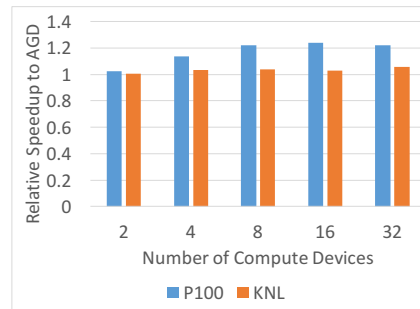


Figure 11: Relative Speedup of Gossip to AGD for P100 and KNL Clusters on CIFAR10 dataset

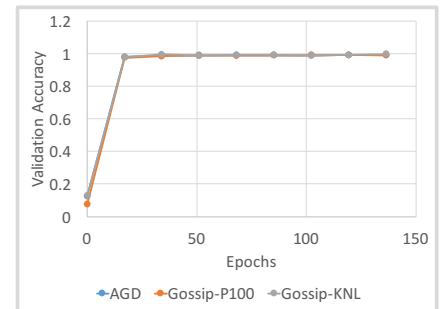


Figure 12: Validation Accuracies of AGD, Gossip on KNL and Gossip on GPU Clusters on MNIST dataset

Name	4	8	16	32	64	128
GossipGraD	100	100	100	100	100	100
PowerAI	100	100	98	99	97	95

Table 7: Compute Efficiency (%) for GossipGraD and PowerAI [13] using up to 128 P100 GPUs. PowerAI performs better than Caffe2, since it uses hierarchical rings based all-to-all reduction

It is worthwhile observing that for GoogLeNet as well, we observe a $\approx 100\%$ computation efficiency at all compute nodes – which implies that GossipGraD is effective for several architectures, and independent of the neural network topology. Figure 16 shows the loss chart for GoogLeNet by comparing AGD and GossipGraD using 32 P100 GPUs. We compare the loss over time. During this time, GossipGraD has only covered 10% of overall iterations – which implies that there is a significant time remaining in the overall training phase. However, even at this short time GossipGraD provides similar or better loss (since lower loss is better) in comparison to AGD.

7.5 Comparing GossipGraD with AGD Communicating Every $\log(p)$ Steps

Another approach to achieve $O(1)$ communication is to combine the models every $\log(p)$ steps instead of every step. Figure 17 compares the performance of this approach to GossipGraD using the LeNet3 network. If the cost of the $\log(p)$ reduction (AGD) can be amortized over $\log(p)$ steps the performance is improved. For the every- $\log(p)$ approach, the performance is trending slightly positive while GossipGraD remains flat. However, the performance of GossipGraD is still greater. Though these two approaches might eventually perform similarly at large scales, the effect on validation accuracy cannot be ignored. Though we expect all approaches to reach target accuracy if tuned with appropriate hyperparameters, for those cases in Figure 17 only GossipGraD was learning. This further shows that GossipGraD is less susceptible to incorrect hyperparameter tuning as one scales.

7.6 Discussion

We evaluated GossipGraD on several dimensions including different datasets, their associated neural network topologies, and GPU/CPU

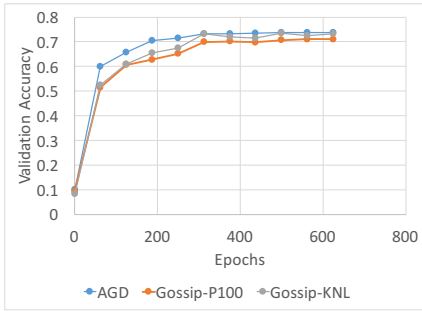


Figure 13: Validation Accuracies of AGD, GossipGrad on KNL and GossipGrad on GPU Clusters for CIFAR10 dataset

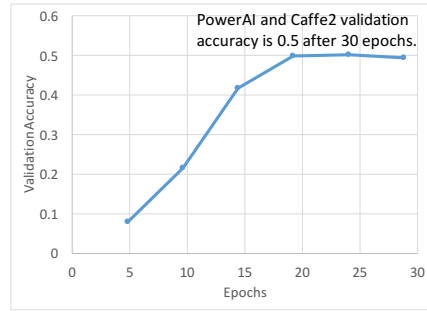


Figure 14: Accuracy of GossipGrad for 128 P100 GPUs for ResNet50

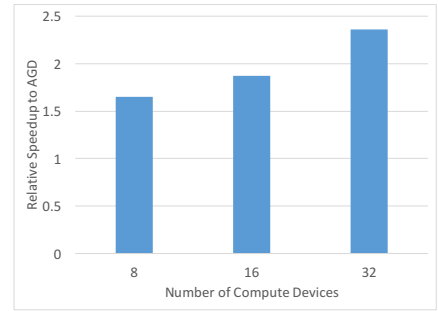


Figure 15: Relative Speedup of GossipGrad to AGD for P100 Cluster for GoogLeNet

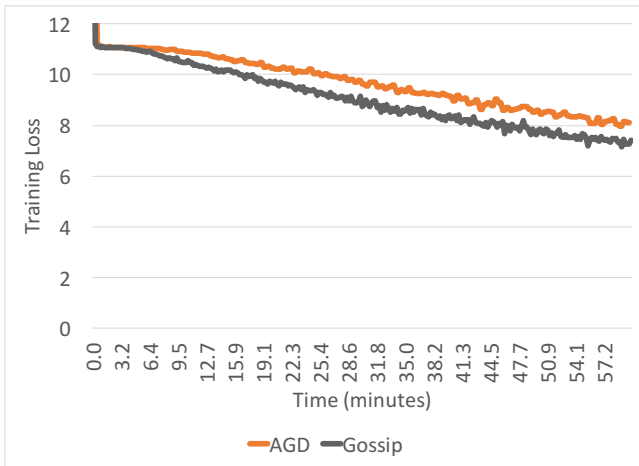


Figure 16: Training Loss after one hour of GossipGrad and AGD on 32 P100 GPUs with GoogLeNet

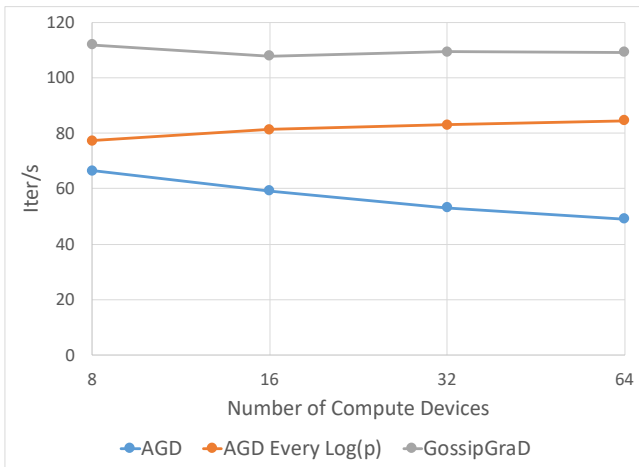


Figure 17: Performance of GossipGrad versus Computing AGD Every Log(p) Iterations

architectures on the metrics of performance and accuracy. We observe that GossipGrad provides complete overlap of communication with computation (as expected), and similar accuracy/validation loss as well studied SGD/AGD. GossipGrad also does not require significant hyper-parameter tuning – which makes it an attractive option for generic datasets. A possible hyper-parameter to tune would be to have a batch size just small enough which can be completely overlapped with computation – however, we have avoided making any changes to batch sizes so that our results are easily reproducible using hyper-parameters available on NVIDIA-Caffe and Intel-Caffe.

8 RELATED WORK

Gradient Descent (GD) is the most widely used algorithm for implementing Deep Learning (DL) algorithms. Since GD is slow, a variant which selects a random subset (*batch*) of the original dataset is used – referred as batch/stochastic gradient descent (SGD). Several frameworks provide high performance implementations of SGD. The most widely used implementations are Caffe [6] (CPUs/GPUs), Warp-CTC (GPUs), Theano [7, 8] (CPUs/GPUs), Torch [9] (CPUs/GPUs), CNTK [57] (GPUs and Distributed Memory using MPI) and Google TensorFlow [5] which use NVIDIA CUDA Deep Neural Network (cuDNN) library. We use SGD as the baseline for designing and implementing GossipGrad and its variants. As evident, deep neural networks (DNN) – which store the model for DL algorithms – suffer from a variety of problems, such as vanishing gradient [58]. Hinton *et al.* [59, 60] proposed a solution to address this problem, by layer-wise training *autoencoders* [61]. Individual layers may be coalesced together to form a single DNN [62].

Several researchers have proposed methods to scale DL algorithms on distributed memory systems. Table 1 shows a table of distributed DL implementations with comparisons on several metrics. We classify these approaches among *parameter server* based and non-parameter server-based approaches. Distbelief [37] is an approach proposed by Dean *et al.*, which uses a parameter server (PS) for model updates at a central server. Several PS optimizations specific to GPU have been proposed in approaches such as GeePS [21] and MXNET [23], and techniques to address delay compensation [63]. Chen *et al.* have studied the limitations of PS and reported that parameter servers quickly become a bottleneck, require expensive *warm up phase* (during which training is conducted on

a single compute node, till loss starts to decrease). However, existing Spark implementations typically use sockets – which is not optimal for HPC interconnects. To leverage HPC interconnects effectively, several Remote Direct Memory Access (RDMA) based PS approaches have been proposed [38, 39]. However, they primarily suffer from typical convergence problems and performance bottlenecks as observed by other researchers [39].

Several researchers have proposed techniques to address the limitations of parameter server based DL implementations [17, 18, 40, 63, 64]. These DL implementations use model/data parallelism. Das *et al.* have proposed hybrid parallelism for scaling DL implementations. However, given the increasing depth of convolutional layers and suggested by Krizhevsky [56], the majority of DL implementations primarily focus on data parallelism – which is the focus of this paper as well.

Data parallelism approaches for scaling DL require hyperparameter tuning as the scale increases [56, 65]. Though this work did not focus on hyperparameter tuning, we expect GossipGrad would also benefit from hyperparameter tuning such as the proposed LARS method [65, 66], RMSprop warm-up [67], and slow-start learning rate schedules [68, 69].

9 CONCLUSIONS

In this paper, we have presented GossipGrad, which is a novel gossip communication protocol based Stochastic Gradient Descent (SGD) algorithm for scaling Deep Learning (DL) algorithms on large-scale systems. GossipGrad has reduced the overall communication complexity from $\Theta(\log(p))$ for p compute nodes in well-studied SGD to $O(1)$ and considered diffusion such that compute nodes exchange their updates (gradients) indirectly after every $\log(p)$ steps. It also considers rotation of communication partners for facilitating direct diffusion of gradients and asynchronous distributed shuffle of samples during the feedforward phase in SGD to prevent over-fitting. We have implemented GossipGrad for GPU and CPU clusters and use NVIDIA GPUs (Pascal P100) connected with InfiniBand, and Intel Knights Landing (KNL) connected with Aries network. We evaluate GossipGrad using well-studied dataset ImageNet-1K ($\approx 250\text{GB}$), and widely studied neural network topologies such as GoogLeNet and ResNet50. Our performance evaluation using both KNL and Pascal GPUs has indicated that GossipGrad can achieve perfect efficiency for these datasets and their associated neural network topologies. Specifically, for ResNet50, GossipGrad is able to achieve $\approx 100\%$ compute efficiency using 128 NVIDIA Pascal P100 GPUs – while matching the top-1 classification accuracy published in literature.

REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in Neural Information Processing Systems 25*, F. Pereira, C. Burges, L. Bottou, and K. Weinberger, Eds. Curran Associates, Inc., 2012, pp. 1097–1105. [Online]. Available: <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>
- [2] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, “Going deeper with convolutions,” in *CVPR 2015*, 2015. [Online]. Available: <http://arxiv.org/abs/1409.4842>
- [3] P. Baldi, P. Sadowski, and D. Whiteson, “Searching for Exotic Particles in High-Energy Physics with Deep Learning,” *Nature Commun.*, vol. 5, p. 4308, 2014.
- [4] Y. Liu, E. Racah, Prabhat, J. Correa, A. Khosrowshahi, D. Lavers, K. Kunkel, M. Wehner, and W. Collins, “Application of deep convolutional neural networks for detecting extreme weather in climate datasets,” *arXiv preprint*

arXiv:1605.01156, 2016.

- [5] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng, “TensorFlow: Large-scale machine learning on heterogeneous systems,” 2015, software available from tensorflow.org. [Online]. Available: <http://tensorflow.org/arXiv> preprint *arXiv:1605.0993*, 2014.
- [6] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, “Caffe: Convolutional architecture for fast feature embedding,” *arXiv preprint arXiv:1408.5093*, 2014.
- [7] J. Bergstra, O. Breuleux, F. Bastien, P. Lamblin, R. Pascanu, G. Desjardins, J. Turian, D. Warde-Farley, and Y. Bengio, “Theano: a CPU and GPU math expression compiler,” in *Proceedings of the Python for Scientific Computing Conference (SciPy)*, Jun. 2010, oral Presentation.
- [8] F. Bastien, P. Lamblin, R. Pascanu, J. Bergstra, I. J. Goodfellow, A. Bergeron, N. Bouchard, and Y. Bengio, “Theano: new features and speed improvements,” Deep Learning and Unsupervised Feature Learning NIPS 2012 Workshop, 2012.
- [9] R. Collobert, S. Bengio, and J. Morio, “Torch: A modular machine learning software library,” 2002.
- [10] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, “Rethinking the inception architecture for computer vision,” *CoRR*, vol. abs/1512.00567, 2015. [Online]. Available: <http://arxiv.org/abs/1512.00567>
- [11] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [12] 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 e-prints*, Jun. 2017.
- [13] M. Cho, U. Finkler, S. Kumar, D. Kung, V. Saxena, and D. Sreedhar, “PowerAI DDL,” *ArXiv e-prints*, Aug. 2017.
- [14] T. Hoefler, T. Schneider, and A. Lumsdaine, “Characterizing the Influence of System Noise on Large-Scale Applications by Simulation,” in *International Conference for High Performance Computing, Networking, Storage and Analysis (SC’10)*, Nov. 2010.
- [15] A. Bhatle, K. Mohror, S. H. Langer, and K. E. Isaacs, “There goes the neighborhood: Performance degradation due to nearby jobs,” in *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*, ser. SC ’13. New York, NY, USA: ACM, 2013, pp. 41:1–41:12. [Online]. Available: <http://doi.acm.org/10.1145/2503210.2503247>
- [16] F. N. Iandola, K. Ashraf, M. W. Moskewicz, and K. Keutzer, “Firecaffe: near-linear acceleration of deep neural network training on compute clusters,” *CoRR*, vol. abs/1511.00175, 2015. [Online]. Available: <http://arxiv.org/abs/1511.00175>
- [17] A. A. Awan, K. Hamidouche, J. M. Hashmi, and D. K. Panda, “S-caffe: Co-designing mpi runtimes and caffe for scalable deep learning on modern gpu clusters,” in *Proceedings of the 22Nd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, ser. PPOPP ’17. New York, NY, USA: ACM, 2017, pp. 193–205. [Online]. Available: <http://doi.acm.org/10.1145/3018743.3018769>
- [18] A. Vishnu, C. Siegel, and J. Daily, “Distributed tensorflow with mpi,” *arXiv preprint arXiv:1603.02339*, 2016.
- [19] H. Zhang, Z. Hu, J. Wei, P. Xie, G. Kim, Q. Ho, and E. P. Xing, “Poseidon: A system architecture for efficient gpu-based deep learning on multiple machines,” *CoRR*, vol. abs/1512.06216, 2015. [Online]. Available: <http://arxiv.org/abs/1512.06216>
- [20] E. P. Xing, Q. Ho, W. Dai, J.-K. Kim, J. Wei, S. Lee, X. Zheng, P. Xie, A. Kumar, and Y. Yu, “Petuum: A new platform for distributed machine learning on big data,” in *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ser. KDD ’15. New York, NY, USA: ACM, 2015, pp. 1335–1344. [Online]. Available: <http://doi.acm.org/10.1145/2783258.2783323>
- [21] H. Cui, H. Zhang, G. R. Ganger, P. B. Gibbons, and E. P. Xing, “Geeps: Scalable deep learning on distributed gpus with a gpu-specialized parameter server,” in *Proceedings of the Eleventh European Conference on Computer Systems*, ser. EuroSys ’16. New York, NY, USA: ACM, 2016, pp. 4:1–4:16. [Online]. Available: <http://doi.acm.org/10.1145/2901318.2901323>
- [22] T. Chilimbi, Y. Suzue, J. Apacible, and K. Kalyanaraman, “Project adam: Building an efficient and scalable deep learning training system,” in *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI 14)*. Broomfield, CO: USENIX Association, 2014, pp. 571–582. [Online]. Available: <https://www.usenix.org/conference/osdi14/technical-sessions/presentation/chilimbi>
- [23] T. Chen, M. Li, Y. Li, M. Lin, N. Wang, M. Wang, T. Xiao, B. Xu, C. Zhang, and Z. Zhang, “Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems,” *CoRR*, vol. abs/1512.01274, 2015. [Online]. Available: <http://arxiv.org/abs/1512.01274>

- [24] “Caffeonspark github project,” <https://github.com/yahoo/CaffeOnSpark>, accessed: 2017-01-23.
- [25] P. Moritz, R. Nishihara, I. Stoica, and M. I. Jordan, “Sparknet: Training deep networks in spark,” *CoRR*, vol. abs/1511.06051, 2015. [Online]. Available: <http://arxiv.org/abs/1511.06051>
- [26] A. Agarwal, E. Akchurin, C. Basoglu, G. Chen, S. Cyphers, J. Droppo, A. Eversole, B. Guenter, M. Hillebrand, R. Hoens, X. Huang, Z. Huang, V. Ivanov, A. Kamenev, P. Kranen, O. Kuchaiev, W. Manousek, A. May, B. Mitra, O. Nano, G. Navarro, A. Orlov, M. Padmilac, H. Parthasarathi, B. Peng, A. Reznichenko, F. Seide, M. L. Seltzer, M. Slaney, A. Stolcke, Y. Wang, H. Wang, K. Yao, D. Yu, Y. Zhang, and G. Zweig, “An introduction to computational networks and the computational network toolkit,” Tech. Rep. MSR-TR-2014-112, August 2014. [Online]. Available: <http://research.microsoft.com/apps/pubs/default.aspx?id=226641>
- [27] P. Chaudhari, C. Baldassi, R. Zecchina, S. Soatto, and A. Talwalkar, “Parle: parallelizing stochastic gradient descent,” *ArXiv e-prints*, Jul. 2017.
- [28] “Paddlepaddle,” <https://github.com/paddlepaddle/paddle>, accessed: 2017-01-23.
- [29] J. Dean, G. S. Corrado, R. Monga, K. Chen, M. Devin, Q. V. Le, M. Z. Mao, M. Ranzato, A. Senior, P. Tucker, K. Yang, and A. Y. Ng, “Large scale distributed deep networks,” in *Proceedings of the 25th International Conference on Neural Information Processing Systems*, ser. NIPS’12. USA: Curran Associates Inc., 2012, pp. 1223–1231. [Online]. Available: <http://dl.acm.org/citation.cfm?id=2999134.2999271>
- [30] “Amazon dsstne github project,” <https://github.com/amzn/amazon-dsstne>, accessed: 2017-01-23.
- [31] S. Tokui, K. Oono, S. Hido, and J. Clayton, “Chainer: a next-generation open source framework for deep learning,” in *Proceedings of Workshop on Machine Learning Systems (LearningSys) in The Twenty-ninth Annual Conference on Neural Information Processing Systems (NIPS)*, 2015. [Online]. Available: http://learningSys.org/papers/LearningSys_2015_paper_33.pdf
- [32] Y. You, A. Buluc, and J. Demmel, “Scaling Deep Learning on GPU and Knights Landing clusters,” *ArXiv e-prints*, Aug. 2017.
- [33] M. Blot, D. Picard, M. Cord, and N. Thome, “Gossip training for deep learning,” *CoRR*, vol. abs/1611.09726, 2016. [Online]. Available: <http://arxiv.org/abs/1611.09726>
- [34] Q. Ho, J. Cipar, H. Cui, J. K. Kim, S. Lee, P. B. Gibbons, G. A. Gibson, G. R. Ganger, and E. P. Xing, “More effective distributed ml via a stale synchronous parallel parameter server,” in *Proceedings of the 26th International Conference on Neural Information Processing Systems*, ser. NIPS’13. USA: Curran Associates Inc., 2013, pp. 1223–1231. [Online]. Available: <http://dl.acm.org/citation.cfm?id=2999611.2999748>
- [35] X. Pan, J. Chen, R. Monga, S. Bengio, and R. Józefowicz, “Revisiting distributed synchronous SGD,” *CoRR*, vol. abs/1702.05800, 2017. [Online]. Available: <http://arxiv.org/abs/1702.05800>
- [36] A. Ichinose, A. Takefusa, H. Nakada, and M. Oguchi, “Pipeline-based processing of the deep learning framework caffe,” in *Proceedings of the 11th International Conference on Ubiquitous Information Management and Communication*, ser. IMCOM ’17. New York, NY, USA: ACM, 2017, pp. 97:1–97:8. [Online]. Available: <http://doi.acm.org/10.1145/3022227.3022323>
- [37] J. Dean, G. Corrado, R. Monga, K. Chen, M. Devin, M. Mao, M. Ranzato, A. Senior, P. Tucker, K. Yang, Q. V. Le, and A. Y. Ng, “Large scale distributed deep networks,” in *Advances in Neural Information Processing Systems 25*, P. Bartlett, F. Pereira, C. Burges, L. Bottou, and K. Weinberger, Eds., 2012, pp. 1232–1240. [Online]. Available: http://books.nips.cc/papers/files/nips25/NIPS2012_0598.pdf
- [38] “Deep image: Scaling up image recognition,” *CoRR*, vol. abs/1501.02876, 2015, withdrawn. [Online]. Available: <http://arxiv.org/abs/1501.02876>
- [39] M. Li, D. G. Andersen, A. Smola, and K. Yu, “Communication efficient distributed machine learning with the parameter server,” in *Proceedings of the 27th International Conference on Neural Information Processing Systems*, ser. NIPS’14. Cambridge, MA, USA: MIT Press, 2014, pp. 19–27. [Online]. Available: <http://dl.acm.org/citation.cfm?id=2968826.2968829>
- [40] D. Das, S. Avancha, D. Mudigere, K. Vaidyanathan, S. Sridharan, D. D. Kalamkar, B. Kaul, and P. Dubey, “Distributed deep learning using synchronous stochastic gradient descent,” *CoRR*, vol. abs/1602.06709, 2016. [Online]. Available: <http://arxiv.org/abs/1602.06709>
- [41] P. H. Jin, Q. Yuan, F. N. Iandola, and K. Keutzer, “How to scale distributed deep learning?” *CoRR*, vol. abs/1611.04581, 2016. [Online]. Available: <http://arxiv.org/abs/1611.04581>
- [42] W. Gropp, E. Lusk, N. Doss, and A. Skjellum, “A High-Performance, Portable Implementation of the MPI Message Passing Interface Standard,” *Parallel Computing*, vol. 22, no. 6, pp. 789–828, 1996.
- [43] A. Geist, W. Gropp, S. Huss-Lederman, A. Lumsdaine, E. L. Lusk, W. Saphir, T. Skjellum, and M. Snir, “MPI-2: Extending the message-passing interface,” in *Euro-Par, Vol. 1*, 1996, pp. 128–135.
- [44] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, “ImageNet Large Scale Visual Recognition Challenge,” *International Journal of Computer Vision (IJCV)*, vol. 115, no. 3, pp. 211–252, 2015.
- [45] S. Marsland, *Machine learning: an algorithmic perspective*. CRC press, 2015.
- [46] S. Sur, H.-W. Jin, L. Chai, and D. K. Panda, “Rdma read based rendezvous protocol for mpi over infiniband: design alternatives and benefits,” in *POPP*, 2006, pp. 32–39.
- [47] Y. N. Dauphin, R. Pascanu, C. Gulcehre, K. Cho, S. Ganguli, and Y. Bengio, “Identifying and attacking the saddle point problem in high-dimensional non-convex optimization,” in *Advances in neural information processing systems*, 2014, pp. 2933–2941.
- [48] H. Robbins and D. Siegmund, “A convergence theorem for non negative almost supermartingales and some applications,” in *Herbert Robbins Selected Papers*. Springer, 1985, pp. 111–135.
- [49] D. Saad, “Online algorithms and stochastic approximations,” *Online Learning*, vol. 5, 1998.
- [50] K. C. Kiwiel, “Convergence and efficiency of subgradient methods for quasi-convex minimization,” *Mathematical programming*, vol. 90, no. 1, pp. 1–25, 2001.
- [51] NVIDIA Corporation, “NVIDIA/caffe, forked from BVLC/caffe,” <https://github.com/NVIDIA/caffe>, 2016.
- [52] Intel Corporation, “intel/caffe, forked from BVLC/caffe,” <https://github.com/intel/caffe>, 2016.
- [53] Y. LeCun, C. Cortes, and C. J. Burges, “The mnist database of handwritten digits, 1998,” *Available electronically at http://yann.lecun.com/exdb/mnist*, 2012.
- [54] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [55] A. Krizhevsky, “Learning multiple layers of features from tiny images,” University of Toronto, Tech. Rep., 2009.
- [56] —, “One weird trick for parallelizing convolutional neural networks,” *CoRR*, vol. abs/1404.5997, 2014. [Online]. Available: <http://arxiv.org/abs/1404.5997>
- [57] A. Agarwal, E. Akchurin, C. Basoglu, G. Chen, S. Cyphers, J. Droppo, A. Eversole, B. Guenter, M. Hillebrand, R. Hoens, X. Huang, Z. Huang, V. Ivanov, A. Kamenev, P. Kranen, O. Kuchaiev, W. Manousek, A. May, B. Mitra, O. Nano, G. Navarro, A. Orlov, M. Padmilac, H. Parthasarathi, B. Peng, A. Reznichenko, F. Seide, M. L. Seltzer, M. Slaney, A. Stolcke, Y. Wang, H. Wang, K. Yao, D. Yu, Y. Zhang, and G. Zweig, “An introduction to computational networks and the computational network toolkit,” Tech. Rep. MSR-TR-2014-112, August 2014. [Online]. Available: <http://research.microsoft.com/apps/pubs/default.aspx?id=226641>
- [58] M. Bianchini and F. Scarselli, “On the complexity of neural network classifiers: A comparison between shallow and deep architectures,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 25, no. 8, pp. 1553–1565, 2014.
- [59] G. E. Hinton and S. Osindero, “A fast learning algorithm for deep belief nets,” *Neural Computation*, vol. 18, p. 2006, 2006.
- [60] Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle, “Greedy layer-wise training of deep networks,” in *Advances in Neural Information Processing Systems 19*, B. Schölkopf, J. C. Platt, and T. Hoffman, Eds. MIT Press, 2007, pp. 153–160. [Online]. Available: <http://papers.nips.cc/paper/3048-greedy-layer-wise-training-of-deep-networks.pdf>
- [61] G. E. Hinton and R. R. Salakhutdinov, “Reducing the dimensionality of data with neural networks,” *Science*, vol. 313, no. 5786, pp. 504–507, Jul. 2006. [Online]. Available: <http://www.ncbi.nlm.nih.gov/sites/entrez?db=pubmed&uid=16873662&cmd=showdetailview&indexed=google>
- [62] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P.-A. Manzagol, “Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion,” *J. Mach. Learn. Res.*, vol. 11, pp. 3371–3408, Dec. 2010. [Online]. Available: <http://dl.acm.org/citation.cfm?id=1756006.1953039>
- [63] S. Zheng, Q. Meng, T. Wang, W. Chen, N. Yu, Z. Ma, and T. Liu, “Asynchronous stochastic gradient descent with delay compensation for distributed deep learning,” *CoRR*, vol. abs/1609.08326, 2016. [Online]. Available: <http://arxiv.org/abs/1609.08326>
- [64] S. De, G. Taylor, and T. Goldstein, “Scaling up distributed stochastic gradient descent using variance reduction,” *CoRR*, vol. abs/1512.02970, 2015. [Online]. Available: <http://arxiv.org/abs/1512.02970>
- [65] Y. You, Z. Zhang, C. Hsieh, J. Demmel, and K. Keutzer, “Imagenet training in minutes,” *CoRR*, abs/1709.05011, 2017.
- [66] Y. You, I. Gitman, and B. Ginsburg, “Scaling sgd batch size to 32k for imagenet training,” *arXiv preprint arXiv:1708.03888*, 2017.
- [67] T. Tieleman and G. Hinton, “Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude,” *COURSERA: Neural networks for machine learning*, vol. 4, no. 2, pp. 26–31, 2012.
- [68] T. Akiba, S. Suzuki, and K. Fukuda, “Extremely large minibatch sgd: Training resnet-50 on imagenet in 15 minutes,” *arXiv preprint arXiv:1711.04325*, 2017.
- [69] 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.

A ARTIFACT DESCRIPTION: [GOSSIPGRAD: GOSSIP PROTOCOL BASED ASYNCHRONOUS GRADIENT DESCENT FOR SCALABLE DEEP LEARNING]

A.1 Abstract

This artifact description contains information needed to reproduce the software environments used by the experiments of the associated submission. The source code contains changes to open source software packages Intel-Caffe and NVIDIA-Caffe.

A.2 Description

A.2.1 Check-list (artifact meta information).

- **Algorithm:** Deep Learning Networks AlexNet and GoogLeNet
- **Program:** Intel-Caffe, NVIDIA-Caffe
- **Compilation:** Sandy-Bridge – intel/16.1.150; Pascal – IBM XL C/C++ for Linux, V13.1.5
- **Data set:** ImageNet ILSVRC2012
- **Hardware:** 2 10-core Sandy-Bridge sockets and 128GB memory per node; 2 10-core IBM POWER8 CPUs 256GB RAM with 4 NVIDIA Tesla P100 GPUs 16GB HBM2 memory per node
- **Experiment workflow:** Clone architecture-specific versions of Caffe, install all caffe dependencies, apply custom patches to add parallel netcdf reader; parallel algorithms; and command-line arguments for caffe tool, compile, run caffe tool with updated solver parameters.
- **Experiment customization:** Upstream-compatible patches were created for Intel-Caffe and NVIDIA-Caffe that add a parallel netcdf DataLayer as well as additional parallel modules. The parallel modules wrap a standard Caffe::Solver instance and implement the Caffe::Solver::Callback interface in order to communicate network data and biases as needed. Lastly, mini-batch sizes were set to 16 instead of their customary defaults for AlexNet and GoogLeNet.
- **Publicly available?:** Yes

A.2.2 How software can be obtained (if available). Intel-Caffe is available from <https://github.com/intel/caffe>. Our software branches from the master branch at commit ID dated 29 December 2016

0bc848bd32d17d5f17bfc7e20915e9805fd1f180 and is available at <https://github.com/matex-org/caffe-intel>.

NVIDIA-Caffe is available from <https://github.com/NVIDIA/caffe>. Our software branches from the master branch at commit ID dated 27 December 2016 6d723362f0f7fe1aaba7913ebe51cc59b12c0634 and is available at <https://github.com/matex-org/caffe-nvidia>.

A.2.3 Hardware dependencies. Intel-Caffe is optimized for Intel CPUs, specifically those featuring advanced vector intrinsics such as AVX. Intel-Caffe automatically downloads a custom MKL library for use with Intel-Caffe if a sufficient version is not automatically located.

NVIDIA-Caffe is specifically tuned for GPUs featuring the latest CUDA and cuDNN versions, specifically CUDA 8.0 and cuDNN 5.1 at the time of manuscript preparation. The Pascal cluster evaluation took advantage of coherent memory between the GPU and CPU for MPI communication using IBM Spectrum MPI, however this feature is not strictly required for the algorithms and advances presented.

A.2.4 Software dependencies. Both Intel-Caffe and NVIDIA-Caffe share the same list of software dependencies with their upstream BVLC-Caffe implementation. IBM also provides a version of Caffe that will compile on the POWER architecture but is otherwise identical to BVLC-Caffe. We copied the specific POWER assembly code for ‘pause’ that was necessary to compile on the Pascal cluster from their github repository at <https://github.com/ibmsoc/caffe>.

A.2.5 Datasets. To evaluate AlexNet and GoogLeNet networks we use ImageNet Large Scale Visual Recognition Challenge 2012 (ILSVRC2012). For downloading ImageNet database please follow url: <http://caffe.berkeleyvision.org/gathered/examples/imagenet.html>. Standard Caffe installations have instructions for creating LMDB databases from these datasets.

We customize the datasets by converting the LMDB databases into netCDF files using a Python script and the pupynere Python package for creating netCDF files. The images were stored as a byte datatype with a corresponding integer for the classification labels. Parallel netCDF was then used for parallel reading of the datasets.

A.3 Installation

Intel-Caffe and NVIDIA-Caffe, with or without our modifications, install as if following the well-known instructions for the upstream BVLC-Caffe with only a few minor exceptions. An MPI compiler is used by setting the CUSTOM.CXX variable in the Makefile.config to use, e.g., mpicxx. The parallel netCDF library is not customarily used with Caffe and must be added to the linker flags. Intel-Caffe recognizes and uses the preprocessor symbol USE.MPI and it must be set to 1.

A.4 Experiment workflow

Caffe uses two related prototxt files per evaluation, one to describe the solver and the other to describe the network. We used prototxt files that come standard with any Caffe distribution. We modified the network prototxt files to use our custom parallel netCDF reader. Otherwise, we modified the network prototxt files to use a batch size of 16 for the majority of our experiments.

The majority of our experiments were evaluating weak scaling and thus used the same prototxt input files. The only exception was when weak scaling the AGD runs. We increased the learning rate each time we doubled the number of compute devices, e.g., Sandy-Bridge nodes, GPUs, based on the observations made by Krizhevsky et al., increasing by a factor of $\sqrt{2}$ each time.

A.5 Evaluation and expected result

Intel-Caffe and NVIDIA-Caffe output is unchanged with respect to BVLC-Caffe. Any traditional analysis of such output remains valid. Any reported accuracy or loss numbers come directly from this standard Caffe output.

A.6 Experiment customization

Besides the changes made to customary batch sizes, the experimentation could be considered as completely customized due to our use of proprietary communication code additions to available open source software packages.