

Theano-MPI: a Theano-based Distributed Training Framework

He Ma¹, Fei Mao², and Graham W. Taylor¹

¹ School of Engineering, University of Guelph, CA {hma02,gwtaylor}@uoguelph.ca

² SHARCNET, Compute Canada, CA feimao@sharcnet.ca

Abstract. We develop a scalable and extendable training framework that can utilize GPUs across nodes in a cluster and accelerate the training of deep learning models based on data parallelism. Both synchronous and asynchronous training are implemented in our framework, where parameter exchange among GPUs is based on CUDA-aware MPI. In this report, we analyze the convergence and capability of the framework to reduce training time when scaling the synchronous training of AlexNet and GoogLeNet from 2 GPUs to 8 GPUs. In addition, we explore novel ways to reduce the communication overhead caused by exchanging parameters. Finally, we release the framework as open-source for further research on distributed deep learning³.

1 Introduction

With the constant improvement of hardware and discovery of new architectures, algorithms, and applications, deep learning is gaining popularity in both academia and industry. Object recognition [20], is now dominated by deep learning methods, which in many cases, rival human performance. Recent success in areas such as activity recognition from video [13] and statistical machine translation [14] is an example of deep learning’s ascent both in performance and at scale.

With the new generations of GPU cards and increased device memory, researchers are able to design and train models with more than 140 million parameters (c.f. VGGNet [21]) and models that are as deep as 150 layers (c.f. ResNet [9]).

The emergence of larger datasets, e.g. ImageNet [20] and MS-COCO [18], challenges artificial intelligence research and leads us to design deeper and more expressive models so that the complexity of models is sufficient for the task.

Despite of the increased computing power of GPUs, it usually takes weeks to train such large models to desired accuracy on a single GPU. This is due to the increased time associated with training deeper models and iterating over the examples in larger datasets. This is where distributed training of deep learning models becomes crucial, especially for activities such as model search which may involve training and evaluating models thousands of times.

³ <https://github.com/uoguelph-mlrg/Theano-MPI>

A naïve approach to scaling up is running several copies of the same model in parallel on multiple computing resources (e.g. GPUs), each computing its share of the dataset and averaging their parameters at every iteration. This strategy is called data parallelism, and its efficient implementation is the focus of our work. More sophisticated forms of distributed training, including model parallelism are important but outside the current scope of our framework.

Theano [23] is an open-source Python library for developing complex algorithms via mathematical expressions. It is often used for facilitating machine learning research. Its support for automatic symbolic differentiation and GPU-accelerated computing has made it popular within the deep learning community. Like other deep learning platforms, including Caffe [12], Torch [3], TensorFlow [1] and MXNet [2], Theano uses CUDA as one of its main backends for GPU accelerated computation. Since a single GPU is limited by its device memory and available threads when solving compute-intensive problems, very recently researchers have started to build multi-GPU support into the most popular frameworks. This includes the multi-GPU version of Caffe (FireCaffe [11]), Torch and Theano (Platoon).

Because the Theano environment usually compiles models for one GPU per process, we need to drive multiple GPUs using multiple processes. So finding a way to communicate between processes becomes a fundamental problem within a multi-GPU framework. There are several existing approaches of implementing inter-process communication besides manually programming on sockets, such as Signals, Message Queues, Message Passing, Pipes, Shared Memory, Memory Mapped Files, etc. However, among those approaches, Message Passing is most suitable for collective communication between multiple programs across a cluster because of its well-developed point-to-point and collective protocols. Message Passing Interface (MPI) is a language-independent communication protocol that can undertake the task of inter-process communication across machines. It is a standardized message-passing system designed for programming on large-scale parallel applications.

Parameter transfer is a basic operation in the distributed training of deep learning models. Therefore, the transfer speed between processes severely impacts the overall data throughput speedup⁴. Since the parameters to be transferred are computed on GPUs, a GPU-to-GPU transfer is required. Compared to the basic `transfer()` function in Theano, NVIDIA GPUDirect P2P technology makes this possible by transferring data between GPUs without passing through host memory. Specifically, it enables CUDA devices to perform direct read and write operations on other CUDA host and device memory. In the context of MPI, GPUDirect P2P technology allows a `GPUArray` memory buffer to be transferred in basic point-to-point and collective operations, making MPI “CUDA-Aware”.

Leveraging CUDA-aware MPI, we have developed a scalable training framework that provides multi-node and multi-GPU support to Theano and efficient inter-GPU parameter transfer at the same time. To the best of our knowledge,

⁴ We define data throughput speedup as the change in total time taken to process a certain amount of examples. It includes both training and communication time.

this is to-date the most convenient way to deploy Theano processes on a multi-node multi-GPU cluster.

2 Related Work

The idea of exploiting data parallelism in machine learning has been widely explored in recent years in both asynchronous and synchronous ways. To accelerate the training of a speech recognition model on distributed CPU cores, DownPour, an asynchronous parameter exchanging method [6], was proposed. It was the largest-scale method to-date for distributed training of neural networks. It was later found that controlling the maximum staleness of parameter updates received by the server leads to faster training convergence [10] on problems like topic modeling, matrix factorization and lasso regression compared to a purely asynchronous approach. For accelerating image classification on the CIFAR and ImageNet datasets, an elastic averaging strategy between asynchronous workers and the server was later proposed [25]. This algorithm allows more exploration of local optima than DownPour and alleviates the need for frequent communication between workers and the server.

Krizhevsky proposed his trick on parallelizing the training of AlexNet [16] on multiple GPUs in a synchronous way [15]. This work showed that eight GPU workers training on the same batch size of 128 can give up to $6.25\times$ data throughput speedup and nearly the same convergence as trained on a single GPU when exploiting both model and data parallelism. Notably, the increase in effective batch size⁵ leads to very small changes in the final convergence of AlexNet when the learning rate is scaled properly. Following his work, a Theano-based two-GPU synchronous framework [7] for accelerating the training of AlexNet was proposed, where both weights and momentum are averaged between two GPUs after each iteration. The model converges to the same level as using a single GPU but in less time.

There has been more development on the acceleration of vision-based deep learning in recent years. NVIDIA developed a multi-GPU deep learning framework, DIGITS, which shows $3.5\times$ data throughput speedup when training AlexNet on 4 GPUs. Purine [17] pipelines the propagation of gradients between iterations and overlaps the communication of large weights in fully connected layers with the rest of back-propagation, giving near $12\times$ data throughput speedup when training GoogLeNet [22] on 12 GPUs. Similarly, MXNet [2] also shows a super-linear data throughput speedup on training GoogLeNet under a distributed training setting.

The Platoon project is a multi-GPU extension for Theano, created and maintained by the official Theano team. It currently supports only asynchronous data parallelism *inside one compute node* based on `posix_ipc` shared memory. In comparison, our framework, Theano-MPI, is designed to support GPUs that are distributed over multiple nodes in a cluster, providing convenient process

⁵ effective batch size = batch size \times number of workers

management and faster inter-GPU memory exchanging based on CUDA-aware MPI.

3 Implementation

Our goal is to make the field of distributed deep learning more accessible by developing a scalable training framework with two key components. First is Theano as a means of constructing an architecture and optimizing it by Stochastic Gradient Descent (SGD). Second is Message Passing Interface (MPI) as an inter-process parameter exchanger. We also aim to explore various ways to reduce communication overhead in parallel SGD and expose some phenomena that affect convergence and speedup when training deep learning models in a distributed framework.

3.1 The BSP Structure

Bulk Synchronous Parallel (BSP) [24] is an intuitive way to implement parallel computing. In the BSP paradigm, workers proceed with training in a synchronous way. Figure 1a shows a 4 GPU example of the proposed BSP structure where the same model is built and run within four processes, P_0, P_1, P_2, P_3 . Each process uses one CPU and one GPU. After the model’s training graph is compiled on the GPU, those parameters in the graph become arrays in GPU memory whose values can be retrieved from device to host and set from host to device. When training starts, the training dataset is split into four parts. In every iteration, each worker process takes a mini-batch of examples from its share and performs SGD on it. After that, all workers are synchronized and model parameters are exchanged between worker processes in a collective way.

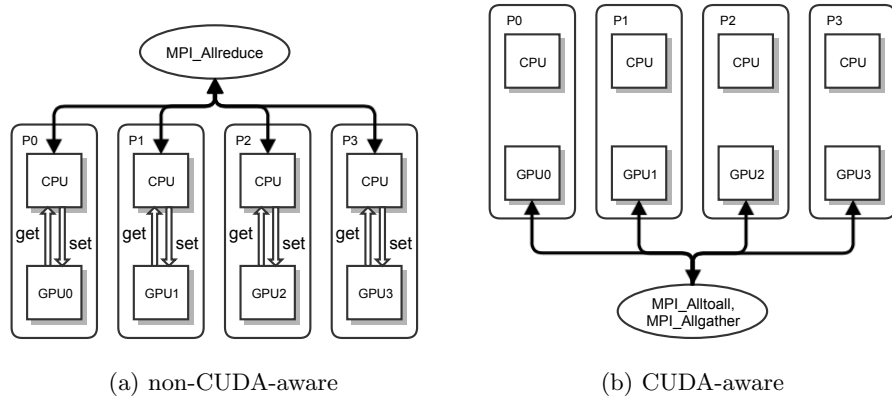


Fig. 1. A 4-GPU example of the BSP structure where arrows indicate communication for parameter exchange.

3.2 CUDA-aware Parameter Exchanging

Synchronous parameter exchange is an array reduction problem which consists of both data transfer and calculation. The GPUDirect P2P technology allows exchanging parameters between GPUs without passing through host memory, making MPI functions “CUDA-aware”. Based on this, we explored various strategies trying to minimize the data transfer and calculation time, and make more efficient use of QPI, PCIe and network card bandwidth during data transfer. The basic strategy is to use the MPI `Allreduce()` function. However, the CUDA-aware version of it in OpenMPI 1.8.7 does not give much improvement since any collective MPI function with arithmetic operations still needs to copy data to host memory. Functions like `Alltoall()` and `Allgather()` do not involve any arithmetic and therefore the CUDA-aware version of them (Fig. 1b) can avoid passing through host memory unless data transfer crossing the QPI bus is needed. We therefore implemented a CUDA-aware `Alltoall-sum-Allgather` strategy which separates the data transfer and computation. An example of this strategy is demonstrated in Fig. 2. Here, the summation kernels required for parameter exchange are executed in parallel on GPUs. Our test shows the GPU summation kernel takes only 1.6% of the total communication time.

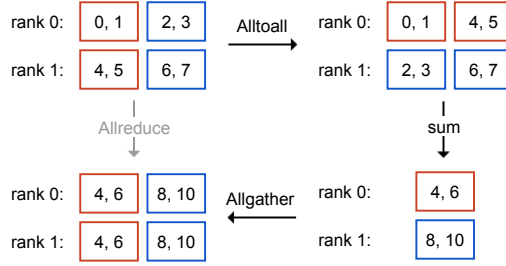


Fig. 2. An example demonstrating the reduction of arrays on rank 0 and rank 1 with the proposed `Alltoall-sum-Allgather` strategy compared to MPI `Allreduce`. Sub-arrays of data items (indicated by same-coloured boxes) need to be summed and the results exchanged with the other ranks.

Using low precision data types for weights or activations (or both) in the forward pass during training of deep neural networks has received much recent interest [4,5]. It was shown that training Maxout [8] networks at 10 bits fixed point precision can still yield near state-of-art test accuracy [4]. In light of this, we also implemented the transfer of parameters at half-precision while summing them at full precision, in order to further reduce communication overhead.

Figure 3 shows the improvement of the combination of strategies over MPI `Allreduce`. The “ASA” strategy shows three times faster communication relative to MPI `Allreduce` and the half precision version of it gives nearly 6 times faster

performance. Those results are obtained on distributed GPUs on 8 nodes in a cluster. Each node hosts one GPU.

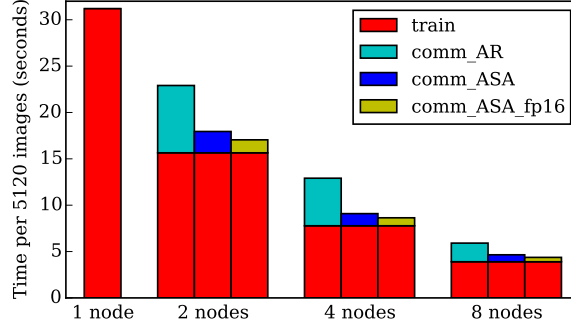


Fig. 3. Computation (train) vs. relative communication overhead of different parameter exchanging strategies during training AlexNet-128b (AR: Allreduce, ASA: CUDA-aware Alltoall-sum-Allgather).

Due to the limitation imposed by the Global Interpreter Lock (GIL) in Python, overlapping the communication with the gradient calculation as in [17] has not yet been implemented in our framework. We expect this, if implemented, would substantially reduce the communication cost of exchanging large matrices in fully-connected layers.

3.3 Parallel Loading

For large-scale visual recognition applications such as ImageNet LSVRC, the data required for training is on the order of hundreds of Gigabytes. Therefore, it is difficult to load all image data completely into memory after training starts. Instead, images are stored as batch files on local or remote disks and loaded one file at a time by each process. Loading image batches x from disk can be time consuming⁶. It is affected by various factors, including file size, file format, disk I/O capability and network bandwidth if reading from remote disks. If in every iteration, the training process should wait for data loading to be ready in order to proceed, one can imagine the time cost by loading data will be critical to the total performance. One way to circumvent this, given the independence of loading and training, is to load those files in parallel with the forward and backward propagations on the last loaded batch. However, this assumes loading one batch of images takes shorter than one iteration of training the model. This auxiliary loading process should follow procedures in Alg. 1 to collaborate efficiently with its corresponding training process:

⁶ Loading labels y , on the other hand, is much faster, therefore labels can be loaded completely into memory.

Algorithm 1 The parallel loading process**Require:**

Host memory allocated for loading image batch $hostdata_x$.
 GPU memory allocated for preprocessed image batch $gpudata_x$
 GPU memory allocated for the actual model graph input $input_x$,
 $mode=$ None, $recv=$ None, $filename=$ None.
 Mean image $image_mean$

Ensure:

```

1: while True do
2:   Receive the  $mode$  (train, validate or stop) from training process
3:   if  $recv=$ "stop" then
4:     break
5:   else
6:      $mode \leftarrow recv$ 
7:   Receive the first filename to be loaded from training process  $filename \leftarrow recv$ 
8:   while True do
9:     Load file "filename" from disk into host memory  $hostdata_x$ .
10:     $hostdata_x = hostdata_x - image\_mean$ 
11:    Crop and mirror  $hostdata_x$  according to  $mode$ .
12:    Transfer  $hostdata_x$  from host to GPU device memory  $gpudata_x$ .
13:    Wait for training on the last  $input_x$  to finish by receiving the next filename
    to be loaded.
14:    if  $recv$  in ["stop", "train", "val"] then
15:      break
16:    else
17:       $filename \leftarrow recv$ 
18:    Transfer  $gpudata_x$  to  $input_x$ .
19:    Synchronize GPU context.
20:    Notify training process to precede with the newly loaded  $input_x$ 

```

Different from the `multiprocessing` and `Queue` messaging method in [7], we used the MPI `Spawn` function to start a child process from each training process and used the resulting MPI intra-communicator to pass messages between the training process and its child process. As shown in Algorithm 1, the parallel loading process can read image files, subtract the mean image, crop sub-images and finally load preprocessed data onto GPUs. By doing this, we are able to overlap the most compute-intensive part (Step 10 to 13 in Algorithm 1) with forward and backward graph propagation in the training process.

4 Benchmarking

Exchanging parameters is a necessary aspect of parallel SGD, however, it can be achieved in a variety of different ways. Parameters updated during SGD include weights (and biases), momentum (if using momentum SGD) and raw gradients. Averaging *weights* after gradient descent (AWAGD) [15,7] is a straightforward parallel SGD scheme. We have proved [19] that training a perceptron using this

scheme on multiple GPUs can either be equivalent to or a close approximate of sequential SGD training on a single GPU depending on whether or not effective batch size is kept constant. In this scheme, the learning rate is scaled with the number of GPUs used [15], namely k . It can also be shown that this scheme is equivalent to summing up the *parameter updates* from all GPUs before performing gradient descent (SUBGD), which does not require scaling up the learning rate. However, our experiments show that tuning the learning rate is still dependent on k to ensure initial convergence of the model. Table 1 lists the learning rates we used and the convergence we achieved in training AlexNet⁷ and GoogLeNet⁸ at different scales (number of workers).

Recent work has applied low precision to weights and activations during training [4]. In the extreme, binary data types have been considered [5]. This enables efficient operation of the low-precision network both at deployment (test time) and during the forward propagation stage during training. However, gradients used for parameter updates must still be stored at high-precision or learning will fail. Related to this, we see a drop in accuracy due to reduced-precision parameter exchange. The validation top-5 error of the “8GPU-32b” AlexNet in Table 1 increased from 19.9% to 20.3% and that of the “8GPU” GoogLeNet increased from 10.65% to 11.7%. Parameter exchange is an important part of the update stage in a distributed training framework. Therefore, there is a necessary accuracy-speed tradeoff to consider when adopting a low-precision strategy.

Table 1. Trade-off between accuracy and speedup under different hyper parameter settings in training AlexNet and GoogLeNet based on the ASA strategy. The learning rate reported was the best one found empirically for the particular setting (HP: hyper-parameters, LR: learning rate, BS: batch size).

# of workers	AlexNet					GoogLeNet				
	HP		Result			HP		Result		
	LR	BS	Accuracy	Speedup		LR	BS	Accuracy	Speedup	
1GPU	0.01	128	19.8%	1×		0.01	32	10.07%*	1×	
2GPU	0.01	128	19.8%	1.7×		0.007	32	10.20%	1.9×	
4GPU	0.01	128	20.4%	3.4×		0.005	32	10.48%	3.7×	
8GPU	0.005	128	20.7%	6.7×		0.005	32	10.65%	7.2×	
8GPU	0.005	32	19.9%	4.9×				-		
8GPU-fp16	0.005	32	20.3%	5.7×		0.005	32	11.75%	7.3×	

⁷ Top-5 error at epoch 62. The implementation is based on theano_alexnet from uoguelph-mlrg. https://github.com/uoguelph-mlrg/theano_alexnet.

⁸ Top-5 error at epoch 70. BVLC GoogLeNet implementation in Caffe is referenced in building the model. https://github.com/BVLC/caffe/tree/master/models/bvlc_googlenet. * At the time of submission, the GoogLeNet 1GPU test is still ongoing. The top-5 error is taken from [22].

Figures 4 and 5 show the convergence of two models trained with SUBGD and the `Alltoall-sum-Allgather` strategy, in which AlexNet is trained on 1, 2, 4 and 8 GPUs with momentum SGD and 128 batch size on each GPU⁹. Similarly, GoogLeNet is trained on 2, 4 and 8 GPUs with a batch size of 32. We see that as more workers are used, the effective batch size becomes too large and the approximation from parallel SGD to sequential SGD becomes worse. As shown in Table 1, one way to preserve convergence at such a large-scale is to reduce the batch size on each worker so that the effective batch size stays small. This gives the model more potential to explore further at low learning rates, though the accuracy improvement at the beginning is slow. However, using smaller batch sizes means more frequent parameter exchanges between workers, which demands attention toward further reducing the communication overhead.

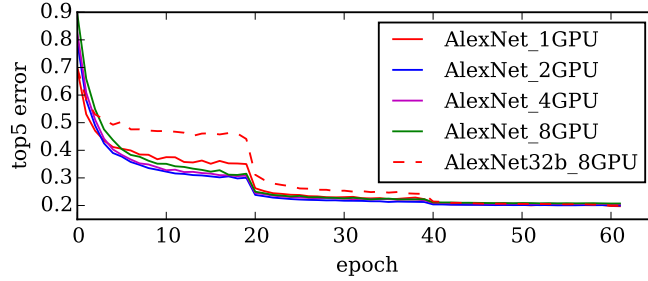


Fig. 4. Validation top-5 error of AlexNet trained at different scales (and batch sizes). Best viewed in colour.

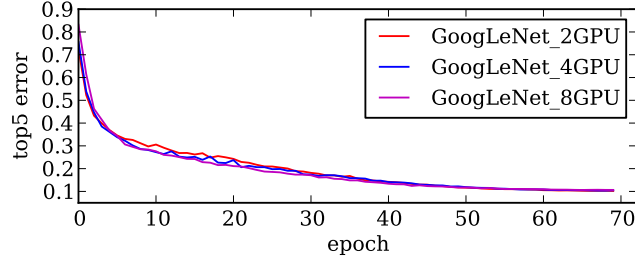


Fig. 5. Validation top-5 error of GoogLeNet trained at different scales. Best viewed in colour.

The speedup of training AlexNet and GoogLeNet are evaluated on 8 distributed GPU nodes (1 GPU per node). To show the performance of accelerating

⁹ Tested on the ILSVRC14 dataset [20].

larger models, we build VGGNet and test its scaling performance on 8 GPUs in a single node with shared memory. This setup meets the memory requirements of VGGNet. Table 2 gives an overview of the structural difference between those three models. Table 3 reports the training and communication time taken to process 5,120 images across different models. We see that these three models scale differently in the framework due to differences in the complexity of their operations as well as the number of free parameters. CUDA-aware parameter exchanging helps boost the speedup of the framework, especially when the number of parameters is relatively large.

Table 2. Structural comparison between the three architectures which were implemented for benchmarking.

Model	Depth ¹⁰	# of parameters ¹¹
AlexNet	8	60,965,224
GoogLeNet	22	13,378,280 ¹²
VGGNet	19	138,357,544

Table 3. Communication overhead per 5,120 images (s) / speedup on 8 GPUs for different models (AR: Allreduce, ASA: CUDA-aware Alltoall-sum-Allgather, ASA16: CUDA-aware Alltoall-sum-Allgather w/ float16).

Model	Train(1GPU)	AR	ASA	ASA16
AlexNet-128b	3.90(31.2)	2.01/5.3×	0.75/6.7×	0.47/7.1×
AlexNet-32b	4.56(36.40)	8.03/2.9×	2.94/4.9×	1.83/5.7×
GoogLeNet-32b	16.82(134.9)	2.07/7.1×	1.96/7.2×	1.76/7.3×
VGGNet-32b	51.79(405.2)	41.41/4.3×	8.60/6.7×	4.84/7.2×

Observing our GoogLeNet ¹³ benchmark result in Fig. 5, we would expect that the framework provides a convergence speedup close to the throughput speedup reported in Table 3, if the convergence of parallel SGD closely ap-

¹⁰ In terms of the amount of parameter-containing layers.

¹¹ In terms of the amount of float32 parameters.

¹² This includes the parameters of the two auxiliary classifiers.

¹³ The learning rate tuning policy of GoogLeNet used was:

$$\eta = \eta_0 \left(1 - \text{EPOCH} \times \frac{\text{NUMBER OF MINIBATCHES}}{\text{MAX ITERATIONS}}\right)^{0.5}.$$

The learning rate tuning policy of AlexNet used was: scaling down by a factor of 10 every 20 epochs.

proximates that of sequential SGD. However, it is difficult to give the exact convergence speedup provided by the framework, since different settings of the hyper-parameters (learning rate tuning policy, weight decay, batch size, cropping randomness) leads to a different convergence path and complicates comparison.

Besides the synchronous framework, we also explored reducing the communication overhead in the asynchronous setting. Referencing the implementation of EASGD in *Platoon*, a Theano-based multi-GPU framework that exploits data parallelism, we re-implemented the framework based on the CUDA-aware MPI `SendRecv()` function without the Round-Robin scheme [25]. Our test shows, when training AlexNet on 8 GPUs, the asynchronous communication overhead in our framework is 42% lower than that in *Platoon* when worker processes communicate with the server in the most frequent way ($\tau = 1$). We also performed a grid search on the hyper-parameters α and τ to achieve better convergence when training AlexNet on eight distributed GPUs, each processing a batch size of 128. The best top-5 error we achieved from this framework was 21.12% at a global epoch of 49 when the moving rate was $\alpha = 0.5$ and averaging period was $\tau = 1$ with a data throughput speedup of $6.7\times$.

5 Hardware and Software Environment

The software was developed and tested on a PI-contributed SHARCNET cluster, named *copper*. As shown in Fig. 6, each node in the cluster is a dual socket system with two NVIDIA Tesla K80 GPUs on each socket. The whole cluster is interconnected with Mellanox Infiniband FDR. We also tested on another cluster, *mosaic*, which features distributed GPUs across nodes connected by Infiniband QDR. Each node has one NVIDIA K20m GPU.

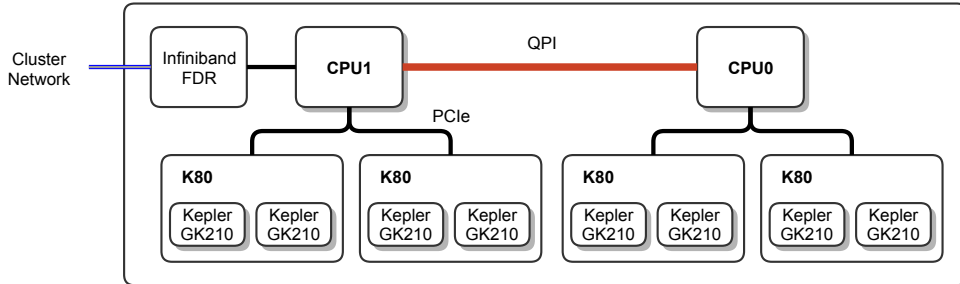


Fig. 6. Hardware connection layout of a copper node

For high-level access to MPI functionality, we use its Python binding `mpi4py`, compiled against OpenMPI 1.8.7. All models mentioned in this report are constructed in Theano 0.8 and their implementation is available in our Github

project. Convolution and pooling operations in the computational graph depend on CUDA 7.0 and the cuDNN v4 library. We also support `cudaconvnet` as an alternative backend.

6 Discussion

We have attempted to scale up the training of deep learning models in an accessible way by developing a scalable training framework built around Theano. Key technical features of our framework are more efficient interprocess communication strategies and parallel data loading techniques. Factors affecting the speedup of the framework can be associated with the model to be trained (i.e. architectural), the training data loading strategy, synchronization in the computational graph, implementation of GPU kernels, system memory and network bandwidth.

Importantly, we try not to compromise the convergence of models trained under our framework since measured speedup is based on the time taken to reach a certain error rate. However, the convergence achieved by a parallel framework also depends on the tuning of that framework’s hyper-parameters. The convergence results in Table 1 can therefore be improved if better hyper-parameters are found. Factors affecting model convergence include the number of worker processes, effective batch size and corresponding learning rate, parameter averaging frequency τ ¹⁴, moving rate α in EASGD and the initialization of model parameters.

The main contributions of our work include: providing multi-node and improved multi-GPU support to the Theano library based on MPI, eliminating substantial communication overhead, exposing convergence and speedup phenomena in parallel SGD, and an implementation of a more efficient parallel loading method.

Our effort towards eliminating the communication overhead involves several aspects: leveraging CUDA-aware MPI for direct data transfer, separating data transfer and summation for more efficient summation on GPUs, and exploring half precision data transfer for faster communication. Our benchmarking results show that our effort on eliminating communication overhead works well on both the 1-GPU-per-node cluster, *mosaic*, and the 8-GPU-per-node cluster, *copper*.

Note that the multi-node testing results in this report are obtained *without* GPUDirect RDMA support due to a limitation in the cluster configuration. Also, the QPI bus topology of a *copper* node limits the usage of GPUDirect P2P technology. This is because the GPUDirect P2P requires all GPUs to be under the same PCIe switch. If a path traversing the QPI is needed, the data transfer would go through CPU RAM first. As a result, further improvement of communication performance based on the current hardware setting would involve consideration of overlapping data transfer with the summation kernel, overlapping parameter exchange with gradient calculation, and designing better

¹⁴ In BSP, we use $\tau = 1$ since larger τ tends to affect convergence in the same way as increasing batch size.

inter-node and intra-node strategies that could balance the bandwidth usage among QPI, PCIe and Infiniband networking.

References

1. Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., et al.: TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems. ArXiv e-prints (Mar 2016)
2. Chen, T., Li, M., Li, Y., Lin, M., Wang, N., et al.: Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems. CoRR abs/1512.01274 (2015)
3. Collobert, R., Kavukcuoglu, K., Farabet, C.: Torch7: A matlab-like environment for machine learning. In: BigLearn, NIPS Workshop (2011)
4. Courbariaux, M., Bengio, Y., David, J.: Low precision arithmetic for deep learning. CoRR abs/1412.7024 (2014)
5. Courbariaux, M., Bengio, Y.: Binarynet: Training deep neural networks with weights and activations constrained to+ 1 or-1. arXiv preprint arXiv:1602.02830 (2016)
6. Dean, J., Corrado, G., Monga, R., Chen, K., Devin, M., et al.: Large scale distributed deep networks. In: Advances in Neural Information Processing Systems 25, pp. 1232–1240 (2012)
7. Ding, W., Wang, R., Mao, F., Taylor, G.W.: Theano-based large-scale visual recognition with multiple gpus. CoRR abs/1412.2302 (2014)
8. Goodfellow, I.J., Warde-Farley, D., Mirza, M., Courville, A., Bengio, Y.: Maxout Networks. ArXiv e-prints (Feb 2013)
9. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. CoRR abs/1512.03385 (2015)
10. Ho, Q., Cipar, J., Cui, H., Lee, S., Kim, J., et al.: More effective distributed ml via a stale synchronous parallel parameter server. In: Advances in Neural Information Processing Systems 26, pp. 1223–1231. Curran Associates, Inc. (2013)
11. Iandola, F.N., Ashraf, K., Moskewicz, M.W., Keutzer, K.: Firecaffe: near-linear acceleration of deep neural network training on compute clusters. CoRR abs/1511.00175 (2015)
12. Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., et al.: Caffe: Convolutional architecture for fast feature embedding. CoRR abs/1408.5093 (2014)
13. Karpathy, A., Toderici, G., Shetty, S., Leung, T., Sukthankar, R., et al.: Large-scale video classification with convolutional neural networks. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2014)
14. Koehn, P., Haddow, B.: Towards effective use of training data in statistical machine translation. In: Proceedings of the Seventh Workshop on Statistical Machine Translation. pp. 317–321. WMT '12, Association for Computational Linguistics, Stroudsburg, PA, USA (2012)
15. Krizhevsky, A.: One weird trick for parallelizing convolutional neural networks. CoRR abs/1404.5997 (2014)
16. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: Advances in Neural Information Processing Systems 25, pp. 1097–1105. Curran Associates, Inc. (2012)
17. Lin, M., Li, S., Luo, X., Yan, S.: Purine: A bi-graph based deep learning framework. CoRR abs/1412.6249 (2014)

18. Lin, T., Maire, M., Belongie, S., Hays, J., Perona, P., et al.: Computer Vision – ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V. Springer International Publishing, Cham (2014)
19. Ma, H.: Developing a Scalable Deep Learning Framework Based on MPI. Master’s thesis, University of Guelph, Guelph, ON, CA (2015)
20. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., et al.: Imagenet large scale visual recognition challenge. *International Journal of Computer Vision* 115(3), 211–252 (2015)
21. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. *CoRR* abs/1409.1556 (2014)
22. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S.E., et al.: Going deeper with convolutions. *CoRR* abs/1409.4842 (2014)
23. Theano Development Team: Theano: A Python framework for fast computation of mathematical expressions. *arXiv e-prints* abs/1605.02688 (May 2016)
24. Valiant, L.G.: A Bridging Model for parallel computation. *Communications of the ACM* 33(8), 103 (1990)
25. Zhang, S., Choromanska, A.E., LeCun, Y.: Deep learning with elastic averaging sgd. In: *Advances in Neural Information Processing Systems* 28, pp. 685–693. Curran Associates, Inc. (2015)