
ACCELERATING NEURAL NETWORK TRAINING WITH DISTRIBUTED ASYNCHRONOUS AND SELECTIVE OPTIMIZATION (DASO)

Daniel Coquelin

Steinbuch Centre for Computing (SCC)
Karlsruhe Institute of Technology (KIT)
Karlsruhe, Germany
daniel.coquelin@kit.edu

Charlotte Debus

Steinbuch Centre for Computing (SCC)
Karlsruhe Institute of Technology (KIT)
Karlsruhe, Germany
charlotte.debus@kit.edu

Markus Götz

Steinbuch Centre for Computing (SCC)
Karlsruhe Institute of Technology (KIT)
Karlsruhe, Germany
markus.goetz@kit.edu

Fabrice von der Lehr

Institute for Software Technology (SC)
German Aerospace Center (DLF)
Cologne, Germany
fabrice.lehr@dlr.de

James Kahn

Steinbuch Centre for Computing (SCC)
Karlsruhe Institute of Technology (KIT)
Karlsruhe, Germany
james.kahn@kit.edu

Martin Siggel

Institute for Software Technology (SC)
German Aerospace Center (DLF)
Cologne, Germany
martin.siggel@dlr.de

Achim Streit

Steinbuch Centre for Computing (SCC)
Karlsruhe Institute of Technology (KIT)
Karlsruhe, Germany
achim.streit@kit.edu

April 16, 2021

ABSTRACT

With increasing data and model complexities, the time required to train neural networks has become prohibitively large. To address the exponential rise in training time, users are turning to data parallel neural networks (DPNN) to utilize large-scale distributed resources on computer clusters. Current DPNN approaches implement the network parameter updates by synchronizing and averaging gradients across all processes with blocking communication operations. This synchronization is the central algorithmic bottleneck. To combat this, we introduce the Distributed Asynchronous and Selective Optimization (DASO) method which leverages multi-GPU compute node architectures to accelerate network training. DASO uses a hierarchical and asynchronous communication scheme comprised of node-local and global networks while adjusting the global synchronization rate during the learning process. We show that DASO yields a reduction in training time of up to 34% on classical and state-of-the-art networks, as compared to other existing data parallel training methods.

Keywords Data parallel neural networks, neural networks, asynchronous communication, hierarchical communication, batch skipping, DASO, MPI, NCCL, optimization, data parallel optimization, HeAT

1 Introduction

Recent advances in deep learning have thrived under the theme "bigger is better". Modern neural networks yield super-human performance on problems such as image classification and semantic segmentation by introducing higher model complexity, for example more layers, inter- and intra-layer connections [1, 2]. However, the training of large networks also requires large datasets. As the sizes of models and datasets increases, so do the computational resources required. In other words, today's deep learning tasks are limited by the hardware and computing time available. In response, parallel training methods have been developed to enable the concurrent use of multiple (distributed) hardware devices.

In general, there are two approaches to parallel training [3]: model parallelism and data parallelism. The model parallel approach distributes the network across multiple computing devices, for example two GPUs with half of the network each. In the data parallel approach, each available computing device trains an identical copy of the network.

Data parallel neural networks (DPNNs) have been used on various architectures and data types to achieve state-of-the-art results [4, 5]. Each model instance in a DPNN performs a forward-backward pass individually over a unique portion of the data, after which the parameters of all networks are synchronized using a global collective operation. This can be effectively viewed as a batch distributed across the devices, i.e. a distributed batch. Traditionally, the synchronization of network parameters is a blocking, averaging operation [3]. This collective blocking operation comprises an inherent bottleneck.

Using non-blocking operations can provide relief as the next forward-backward step can begin while communication is ongoing. However, as global parameter updates are running asynchronously, parameters found by individual network instances are always slightly out-of-date. Out-of-date parameters can also be referred to as stale.

Although computing devices can take many forms, GPUs are currently the most efficient and powerful for training neural networks. Therefore, we will refer to computing devices as GPUs throughout this paper.

Averaging network parameters across multiple instances, traditionally referred to as mini-batch optimization, is only an approximation of the true gradients that would be calculated over the unified batch with batch optimization. Moreover, the standard communication structure communicates with each GPUs individually to synchronize network parameters. This neglects the structure of most computer clusters, where multiple GPUs are grouped on computing nodes with significantly faster node-local connections as compared to cross-node communication.

Large multi-node DPNNs can instead be divided into node-local DPNNs which are themselves members of a global DPNN. This hierarchical approach would significantly reduce the communication overhead, as less data is sent between nodes. Furthermore, what if global parameter synchronization did not occur after every batch and instead the average was calculated asynchronously every B^{th} batch?

To this end, we present our key contribution: the distributed selective and asynchronous optimization (DASO) method. DASO performs communication for network parameter updates in a hierarchical manner: on the node-local level, in the form of GPU-to-GPU operations, and on the global level, where computing nodes are treated as individual entities. This approach allows DASO to perform the time-expensive global synchronization after multiple batches instead of after every forward-backward pass, thus leveraging the potential of acceleration via parallel computation on modern computer clusters.

The remainder of this paper is organized as follows. In Section 2 we will discuss relevant work previously done in the area of data parallel model training. Section 3 introduces the concept of selective distributed asynchronous optimization, followed by performance evaluations on the tasks of image classification and semantic segmentation in Section 4. Our results are summarized and discussed in Section 5, which also gives an outlook towards further improvement and application of the method.

2 Related Work

Data parallel neural networks are the go-to option for accelerating training on large datasets. In DPNNs, each local network is optimized locally, e.g using SGD, before the optimization results are synchronized with all other networks. The most straightforward approach to global synchronization is a collective blocking, average operation after every forward-backward step. This inherently limits the speed of the data parallel training.

Recently, advancements have been made in accelerating the synchronization process by starting the communication of gradient updates while the backward pass is ongoing, with one reporting training times of only 74.7 s on the ImageNet data set [4]. However, this is a tailored approach which does not generalize, as it is highly optimized for a specific network and requires considerable tuning.

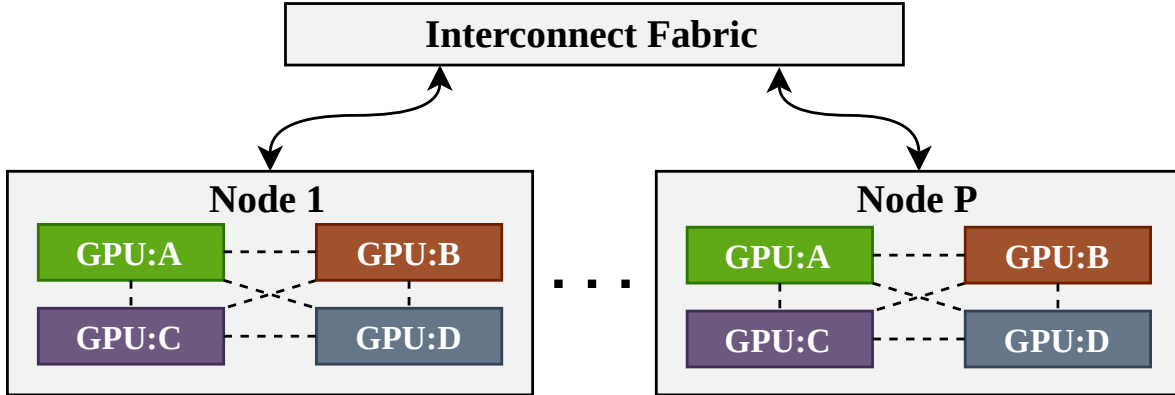


Figure 1: An overview of a common node-based computer cluster with P nodes and four GPUs per node. GPU colors represent *group* membership. The dashed lines indicate GPU-to-GPU communication channels.

Several works have investigated the use of asynchronous SGD (ASGD) [6, 7, 8], which updates the parameters whenever a network finishes a backward pass. Each network retrieves the current model parameters from a parameter server before performing a forward-backward pass. After finishing the backward step, the network sends its updated parameters back to the server, which determines the new global parameters using the updates from all processes. However, if a network is still computing the forward-backward pass when the parameter server is updated, the network’s current parameters are outdated. The subsequently found gradients are referred to as stale. Stale gradients can be leveraged to approximate accurate network parameters, and ASGD has been shown to yield consistent convergence [9]. Recent attempts at accelerating ASGD have been made using individual network optimizers for a warm-up phase and delayed updates to the parameter server [10].

PyTorch [11] and TensorFlow [12] are currently the largest machine learning frameworks. Both offer options for traditional data parallel training. For large systems, a global communication protocol, such as MPI [13], is often required to leverage specialized inter-node connections. Recently, there have been many advancements in the optimization of the global parameter synchronization operation by using MPI with multiple network topologies [14, 15]. These approaches have shown promising results, but remain centered around the idea of a global synchronization for each forward-backward pass.

Currently, the most popular MPI-enabled DPNN framework is Horovod [16]. To reduce the size of data sent via the communication network, Horovod uses tensor fusion, or grouping parameters together to be communicated in a larger chunk of data, and data compression. The data compression in Horovod is frequently done by casting the network parameters into 16-bit floating-point format.

3 Distributed Asynchronous and Selective Optimization (DASO)

The common approach to training DPNNs is to perform a forward-backward pass on each network instance with one portion of the distributed batch, then synchronize the network parameters via a global averaging operation. The averaging of gradients is only an approximation of the true gradients that would be calculated for the entire batch when processes on a single GPU. This approximation is made under the assumption that each portion of the distributed batch is independent and identically distributed (iid) [17].

Under the iid assumption, another approximation can be made: the average parameters of a subset of networks are not significantly different than the average parameters of the complete set of networks. Recalling that modern HPC clusters have different inter- and intra-node communication capabilities (with different bandwidths and latencies), we can utilize this approximation to reduce the communication needed for parallel training, thereby alleviating the intrinsic bottleneck of blocking synchronizations.

We therefore propose the Distributed Asynchronous and Selective Optimization (DASO) method. Instead of a uniform communications network across multiple multi-GPU nodes, DASO uses a hierarchical network model with node-local networks and a global network.

The global network spans all GPUs on all nodes, while the node-local networks are composed of the GPUs on each individual node. The global network is divided into multiple *groups*, with each *group* containing a single GPU from

every node. Global communication takes place exclusively within a *group*, i.e. only *group* members exchange data, while members of other *groups* do not participate. Communication between the node-local GPUs is then handled by the local network, which benefits from high-speed GPU-to-GPU interconnects and optimized communication packages (e.g. NCCL [18]). Under the assumption that the cluster node configurations are homogeneous, DASO creates *groups* between GPUs with the same local identifier as is shown in Figure 1. With this approach, inter-node communication can be reduced by a factor equal to the minimum number of GPUs per node.

Local Synchronization

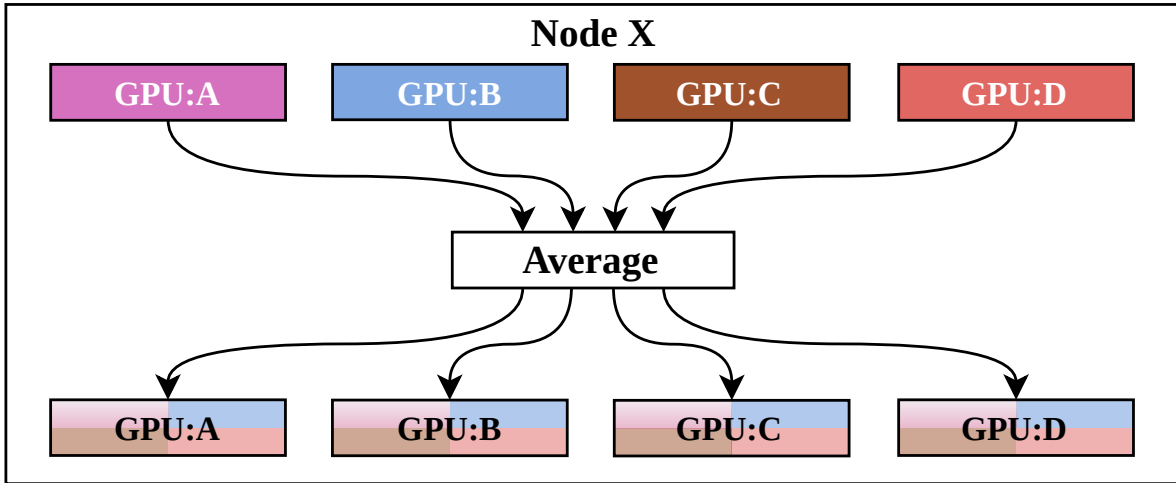


Figure 2: Schematic of the local synchronization step for a single node with four GPUs. The gradients from each GPU are averaged, then each GPU’s gradients are set to the result.

DASO utilizes a multi-step synchronization. Local synchronization (Figure 2) occurs after each batch and uses the node-local network to do gradient-averaging between the local GPUs. Global synchronization (Figure 3) occurs after one or more local synchronizations, in which the network parameters of all members of a single global *group* are shared and averaged. Following every global synchronization, a local update step broadcasts averaged parameters from the local *group* member to all other node-local GPUs (Figure 4). The role of global synchronization rotates between *groups* to overlap communication and computation. Global synchronization can be performed in a blocking or non-blocking manner. In the blocking case, all synchronization steps are performed after each batch. To reduce the amount of data transferred, parameters are cast to a 16-bit datatype representation during buffer packaging. This operation does not effect convergence, as shown by [19]. Once received, the parameters are cast back to their original datatype. In the non-blocking case, the next forward-backward pass is started after the parameters are sent but before they are received. Datatype casting is not beneficial in this scenario, as it delays the start of parameter communications. Each neural network will conduct B forward-backward passes complete with local synchronization before the *group* members receive the sent parameters. Hence, the updates from the global communication step are outdated upon their arrival. To compensate for this, a weighted average of the stale global parameters and the current local parameters is calculated as follows:

$$x_{t+S} = \frac{2Sx_{t+S-1}^l + \sum_{i=1}^P x_t^i}{2S + P} \tag{1}$$

where x_{t+S}^l is the model state on GPU l after S batches to wait after starting batch t for the global synchronization data, x_{t+1}^i is the globally exchanged model states, and P is the number of GPUs in the global network. The weighting of the local parameters was found experimentally. A detailed explanation of Equation (1) and its validity is provided in the supplementary material.

Training of a network with the DASO method can be divided into three key phases: warm-up, cycling, and cool-down. The warm-up and cool-down phases utilize blocking global synchronizations, while the cycling phase uses non-blocking global synchronizations. Given a fixed number of total epochs, warm-up and cool-down phases occur for a set number of epochs at the beginning and end of training respectively. The warm-up phase is used to quickly move away from the randomly initialized parameters and prepare for the cycling phase. The cool-down phase is intended to reduce the slight errors which can arise due to the slight deviation from the iid assumption for individual batches.

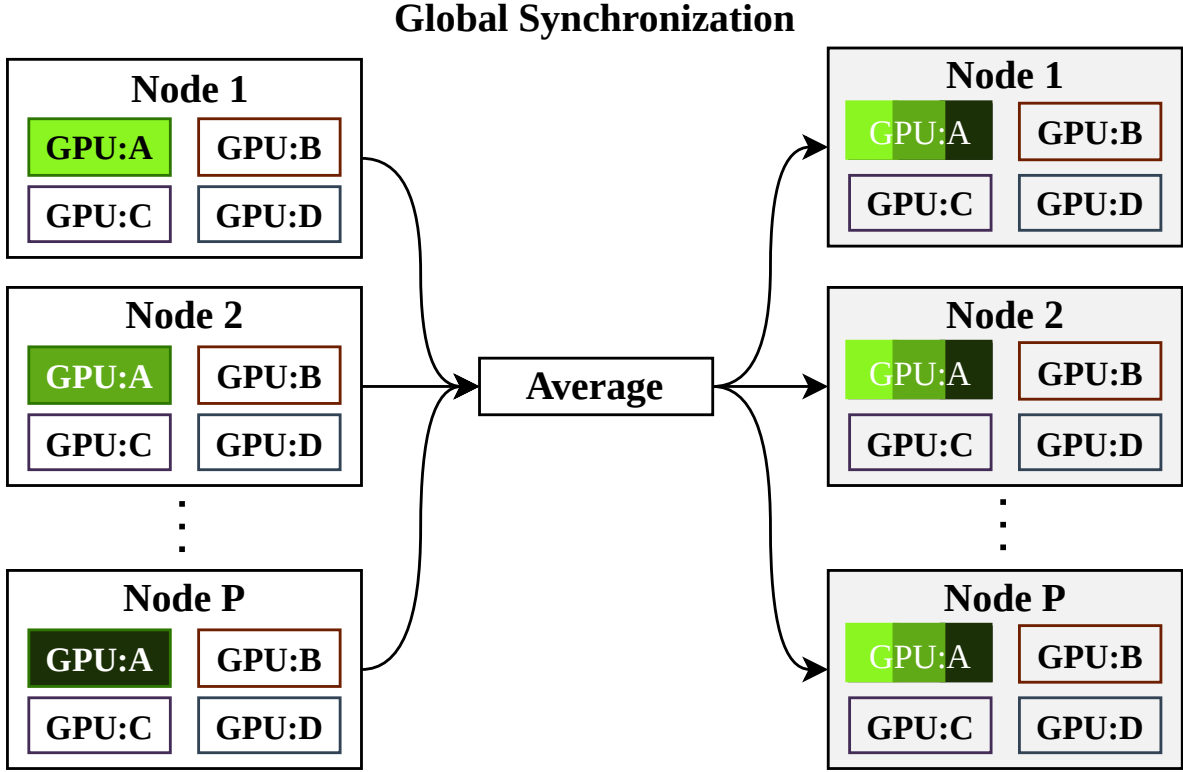


Figure 3: Schematic of the global synchronization step performed by the global communication *group* consisting of GPU:A on each node. The network parameters are averaged by each GPU in the *group*, and the network parameters of each *group* member are set to the result.

In the cycling phase, the number of forward-backward passes between global synchronizations (B) and the number of batches to wait for global synchronization data (W) are varied. B is specified manually upon initialization. For W , an initial value of $B/4$ was found empirically to perform best. Each time the training loss plateaus, B and W are reduced by a factor of two, down to a minimum of one. When $B, W = 1$ and the loss has plateaued, both are reset to their initial values and the process is repeated until the cool-down phase. The synchronization steps in the cycling phase are schematically shown in Figure 5.

3.1 Current Implementation

A DASO proof-of-concept is currently implemented in the HeAT framework [20] for usage with PyTorch networks. HeAT is an open-source Python framework for distributed and GPU-accelerated data analytics which offers both low level array computations as well as assorted higher-level machine learning algorithms. The local networks utilize PyTorch’s DistributedDataParallel class and distributed package [21]. The global communication network utilizes HeAT’s MPI backend, which handles the automatic communication of PyTorch Tensors. The global *groups* are implemented as MPI groups.

To use this implementation of DASO to train an existing PyTorch network, only four additional functions need to be called and the data loaders need to be modified to distribute the data between all GPUs¹. The function calls are illustrated in Listing 1. First, the node-local PyTorch processes are created, which will be utilized during the local synchronization step. Next, the optimizer instance, i.e. DASO, is created with a PyTorch node-local optimizer (e.g. SGD) and the number of epochs for training is specified. The DASO instance will find the aforementioned PyTorch processes automatically.

¹The data loaders need only know how many GPUs exist and what their global rank is.

Local Update

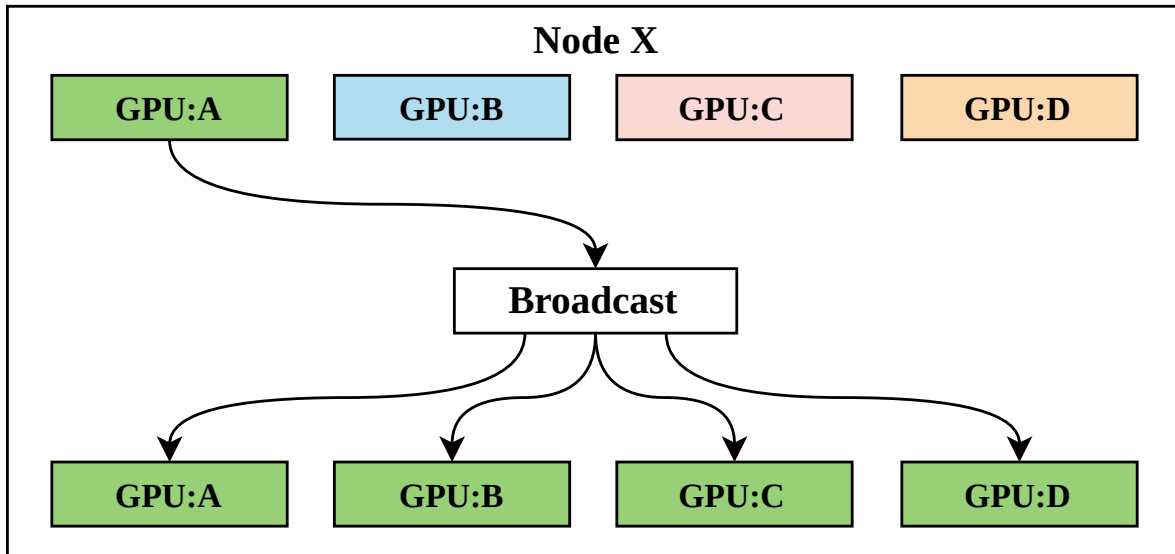


Figure 4: Schematic of the local update step to be performed after the global synchronization step shown in Figure 3. The group member responsible for the global communication, in this case GPU:A, sends its network parameters to all other node-local GPUs, which replace the old parameters on those GPUs.

Listing 1: Simplified training script demonstrating the usage of DASO in HeAT for a PyTorch neural network (net) and PyTorch optimizer (optimizer).

```

1 import heat as ht
2 import torch
3 ...
4 # create PyTorch distributed group
5 world_size = ht.MPI_WORLD.size
6 rank = ht.MPI_WORLD.rank
7 local_rank = rank % num_local_gpus
8 torch.distributed.init_process_group(
9     backend="nccl",
10    rank=local_rank,
11    world_size=world_size
12 )
13 ...
14 # the DASO optimizer is created
15 daso_optimizer = ht.optim.DASO(
16     local_optimizer=optimizer,
17     total_epochs=num_epochs
18 )
19 ...
20 # the hierarchical network is created
21 ht_model = ht.nn.DataParallelMultiGPU(
22     net,
23     daso_optimizer
24 )

```

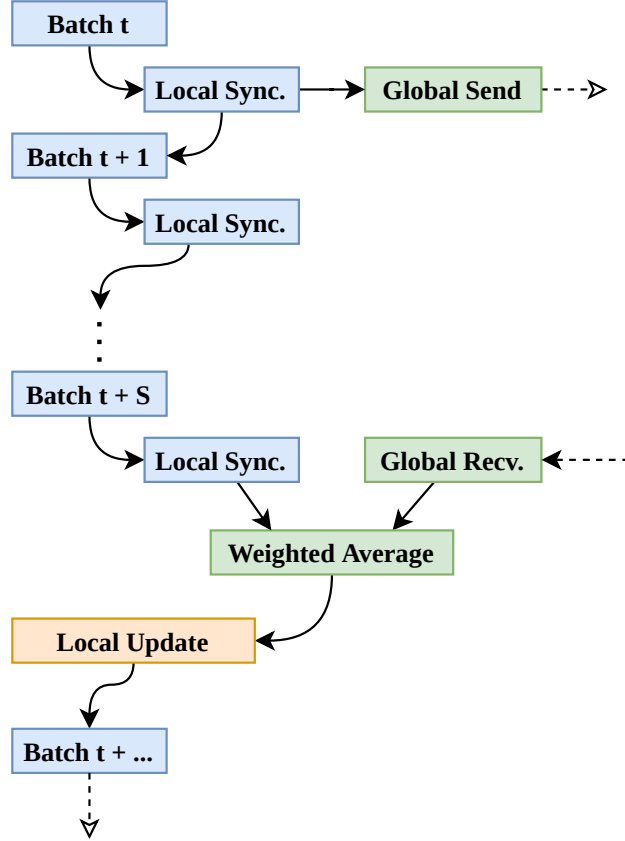


Figure 5: Process flow diagram of the synchronization steps during the cycling phase where t is the batch number and S is the batches to wait before global synchronization. The weighted average is calculated as shown in Equation (1)

4 Performance Evaluation

We evaluate the DASO method on two common examples of data-intensive neural network challenges: a) image classification and b) semantic segmentation. For image classification, we trained ResNet-50 [1] on the ImageNet-2012 [22] dataset. This can be considered a standard benchmark for machine learning, since pre-trained ResNet-50 networks are the backbone of many computer vision pipelines [23]. For semantic segmentation, we trained a state-of-the-art hierarchical multi-scale attention network [5] on the CityScapes [24] dataset.

All experiments were conducted on the JUWELS Booster at the Jülich Supercomputing Center [25]. This center’s HPC cluster has 936 GPU nodes each with two AMD EPYC Rome CPUs and four NVIDIA A100 GPUs, connected via an NVIDIA Mellanox HDR InfiniBand interconnect fabric. The following software versions were used: CUDA 11.0, ParaStationMPI 5.4.7-1-mt, Python 3.8.5, PyTorch 1.7.1+cu110, Horovod 0.21.1, and NCCL 2.8.3-1. The JUWELS Booster provides a CUDA-aware MPI implementation, meaning that GPUs can communicate directly with other GPUs.

We compared DASO to Horovod, as this is currently the most popular choice for MPI-based parallel training of neural networks on computer clusters. We elected not to compare with PyTorch’s distributed package as it utilizes a similar approach to Horovod, namely compression and bucketing. Comparisons are done with respect to training time and accuracy.

Relevant network hyperparameters remain consistent for DASO and Horovod for each experiment. All tested networks use a learning rate scheduler. When the training loss plateaus, i.e. the training loss is not decreasing by more than a set percentage threshold, the scheduler decreases the learning rate by a set factor. Settings of the scheduler, as well as for the local optimizer settings, were set to be identical for both DASO and Horovod for each use-case. With respect to message packaging, Horovod was configured to use floating point 16 compression while DASO compresses to brain floating point 16. The batch size of training is fixed for each GPU in all experiments. Therefore, the combined

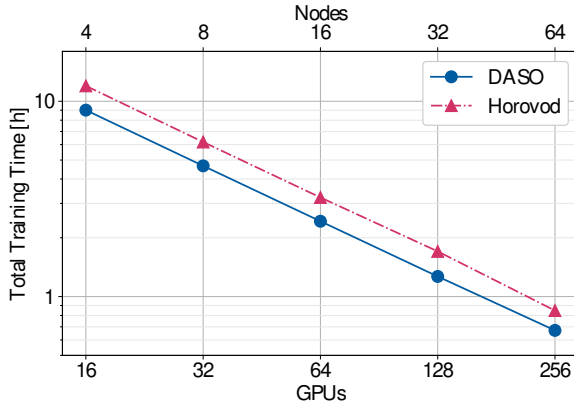


Figure 6: ResNet-50 training time on the ImageNet dataset with DASO and Horovod for increasing node counts. Each node has 4 GPUs.

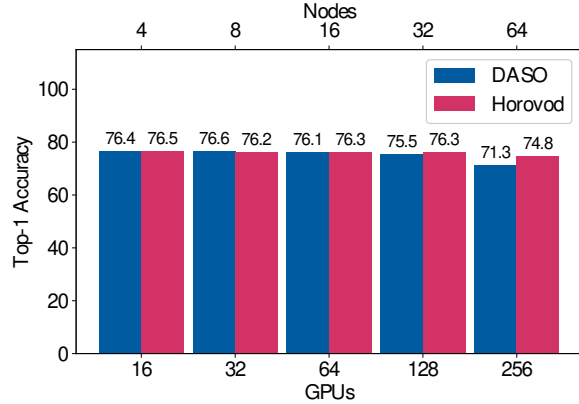


Figure 7: Top-1 accuracy of ResNet-50 networks on the ImageNet dataset when trained with DASO and Horovod for various node counts, each with 4 GPUs.

distributed batch size increases by the number of GPUs times the local batch size. DASO’s maximum number of batches between global synchronizations was set to four for both experiments.

4.1 Image Classification – ImageNet

This experiment was conducted using the ResNet-50 architecture on the ImageNet dataset [22]. For this experiment, the ImageNet-2012 is a large dataset containing 1.2 million labeled images. We evaluate classification quality using top-1 accuracy, i.e. the accuracy with which the model predicts the image labels correctly with a single attempt. For training ResNet-50 on the ImageNet dataset, we consider a 75% top-1 accuracy to be a successful training.

File loading from disk and preprocessing were done using DALI [26]. Training was conducted using cross entropy loss and SGD with a momentum of 0.9 and weight decay of 0.0001 for 90 epochs with a learning rate warm-up phase of five epochs. These values were adapted from PyTorch’s example training script for ResNet-50 on ImageNet. The maximum learning rate is scaled with the number of global processes. The learning rate decays by a factor of 0.5 when the training cross entropy loss is stable for 5 epochs.

Training was conducted on 4, 8, 16, 32, and 64 nodes, which equals 16, 32, 64, 128, and 256 GPUs, respectively. This corresponds to traditional strong scaling experiments for parallel algorithms, where an increase in nodes should ideally result in a proportional reduction in time.

Results of the experiment are shown in Figure 6. Both DASO and Horovod show desirable strong scaling behavior, i.e. a factor of two in GPU number results in the training time being halved. Due to DASO’s optimized hierarchical communication scheme and the reduced number of synchronizations, DASO requires up to 25% less time for training compared to Horovod.

It can further be observed that up to 128 GPUs, DASO and Horovod yield similar levels of accuracy, see Figure 7. However, with more than 128 GPUs, both approaches did not exceed 75% top-1 accuracy. This is due to the fact that accuracy starts to decrease at larger batch sizes in a traditional network unless special allowances are made [27]. Since we keep the portion of the distributed batch that is processed on each individual GPU the same, larger GPU counts ultimately result in a larger distributed batch. Hence, accuracy ultimately decreases. For DASO, the effect is more dramatic because completing batches without a global synchronization has a similar effect to increasing the size of the batch.

4.2 Semantic Segmentation – CityScapes

To further evaluate the performance of the DASO method, we conducted experiments on a cutting edge, state-of-the-art network. To this end, a hierarchical multi-scale attention network [5] was trained for semantic segmentation on the CityScapes [24] dataset. This dataset comprises a collection of images of streets in 50 cities across the world, with 5,000 finely annotated images and 20,000 coarsely annotated images. The network has an HRNet-OCR backbone, a dedicated fully convolutional head, an attention head, and an auxiliary semantic head [5].

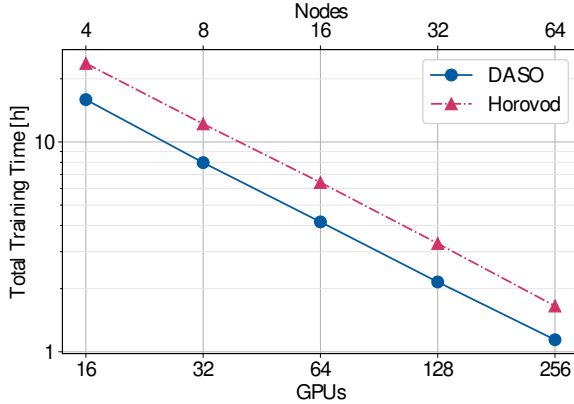


Figure 8: Training time for the selected hierarchical split level attention network [5] on the CityScapes dataset with DASO and Horovod for increasing node counts, each with 4 GPUs.

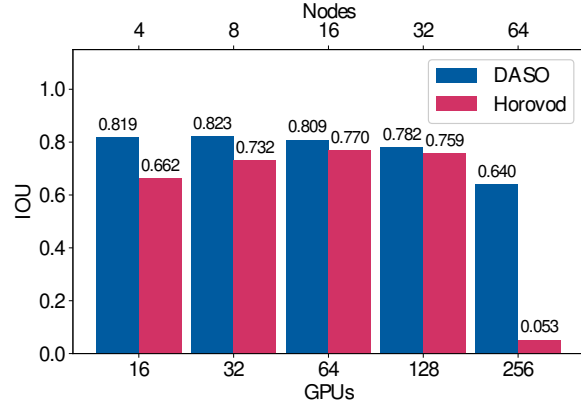


Figure 9: Maximum IOU of the hierarchical split level attention networks on the CityScapes dataset when trained with DASO and Horovod with various node counts, each with 4 GPUs.

The quality of semantic segmentation networks is often evaluated based on the intersection over union (IOU) [28] score. IOU is defined as the intersection of the correctly predicted annotations with the ground truth annotations, divided by their union. In this work, IOU ranges from 0.0 to 1.0.

The network was trained using the following parameters: 175 epochs; the region mutual information loss [29] function; a local SGD optimizer with a weight decay of 0.0001, a momentum of 0.9, and an initial learning rate of 0.0125; a learning rate scheduler which decays the learning rate by a factor of 0.75 when the loss is judged to be stable for 5 epochs. The number of epochs, loss function, and optimizer settings were determined by the original source [5]. The learning rate scheduler deploys a warm up phase of 5 epochs, in which the learning rate is slowly increased from 0.0 to 0.4, after which it decays as scheduled. For the DASO experiments, the synchronized batch normalization operation is conducted within the node-local process group.

In its original publication, the network was trained using supplementary data, whereas the herein presented experiments are performed using only the CityScapes dataset. To determine a baseline accuracy, the original network was trained with four GPUs on a single node using PyTorch’s DistributedDataParallel package. This baseline measurement employed a polynomial decay learning rate scheduler, PyTorch’s automatic mixed precision training and synchronized batch normalization layers. For more detail, see [5]. The baseline IOU of the original network was found to be 0.8258.

During the experiments, we found that for Horovod neither the automatic mixed precision nor the synchronized batch normalization functioned as intended when using the system scheduler software (SLURM [30]). Horovod requires usage of its custom scheduler horovodrun to enable full feature functionality. However, this software is not natively available on many computer clusters, including the JUWELS booster supercomputer. Hence, automatic mixed precision was removed and the synchronized batch normalization layers were replaced with local batch normalization layers.

Training times for various node counts are shown in Figure 8. For up to 128 GPUs, DASO completed the training process in approximately 35% less time than Horovod, demonstrating the advantage of our approach to fully leverage the systems communication architecture together with asynchronous parameter updates. At higher GPU counts the time savings drop to 30%, because there are fewer batches per epoch and hence skipping global synchronization operations provides less benefits.

Quality measurements (IOU) are shown in Figure 9. Although there is a very clear difference between Horovod and DASO, neither matches the accuracy of the baseline network. This is due to the naive learning rate scheduler used for training. With a tuned learning rate optimizer the 16, 32, and 64 node configuration should more accurately recreate the results of the baseline network. At 256 GPUs, training with Horovod did not yield any meaningful results. We hypothesize that this is caused by the lack of a functioning synchronized batch normalization operation in combination with a very large mini-batch.

5 Conclusion

In this work, we have introduced the distributed asynchronous and selective optimization (DASO) method. DASO utilizes a hierarchical communication scheme to fully leverage the communications infrastructure inherent to node-based computer clusters, which often see multiple GPUs per node. By favoring node-local parameter updates, DASO is able to reduce the amount of global communication required for full data parallel network synchronization. Thereby, our approach alleviates the bottleneck of blocking synchronization used in traditional data parallel approaches. We show that, if independent and identically distributed (iid) batches can be reasonably assumed, the global synchronization ubiquitous to the training of DPNNs is not required after each forward-backward pass.

We evaluated DASO on two common DPNN use-cases: image classification on the ImageNet dataset with ResNet-50, and semantic segmentation on the CityScapes dataset with a cutting edge multi-head attention network architecture. Our experiments show that DASO can reduce training time by up to 34% while maintaining similar prediction accuracy when compared to Horovod, the current standard for MPI-based data parallel network training.

At large node counts, DASO and Horovod both suffer a decrease in network accuracy. This is a well-known problem which relates to an increase in the distributed batch size. The effect is more pronounced with DASO due to the reduced number of global synchronization steps. This allows for the identification of where network modifications must be employed to handle very large node counts. We also note that DASO and Horovod will both yield sub-optimal results on datasets for which the iid assumption no longer holds. For those cases, however, data parallel training will be ineffective regardless of the communications scheme. Overall, DASO achieves close-to-optimal accuracies significantly faster than Horovod. Therefore, DASO is optimal for rapid initial training of large networks/datasets, where the training can be further fine-tuned using more traditional methods.

Ultimately, DASO improves the scalability of data parallel neural networks and demonstrates that using more GPUs does not have to be the only solution to speeding up training. With DASO, it is possible to efficiently train large models or process more training data. The beauty of our approach lies in the fact that it is a generic, non-tailored, and easy to implement approach that translates well to any large scale, node-based computer cluster or high-performance computing system. DASO opens the door to redefining data parallel neural network training towards asynchronous, multifaceted optimization approaches.

6 Acknowledgments

This work is supported by the Helmholtz Association Initiative and Networking Fund under project number ZT-I-0003, the Helmholtz AI platform grant and the HAICORE@KIT partition.

References

- [1] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778. IEEE, June 2016.
- [2] Ashish Vaswani, Noam Shazeer, Niki Parmar, et al. Attention Is All You Need, 2017. [accessed on 2021-01-26].
- [3] Tal Ben-Nun and Torsten Hoefler. Demystifying Parallel and Distributed Deep Learning: An In-depth Concurrency Analysis. *ACM Computing Surveys (CSUR)*, 52(4):1–43, 2019.
- [4] Masafumi Yamazaki, Akihiko Kasagi, Akihiro Tabuchi, Takumi Honda, et al. Yet Another Accelerated SGD: ResNet-50 Training on ImageNet in 74.7 seconds, 2019. [accessed on 2021-01-26].
- [5] Andrew Tao, Karan Sapra, and Bryan Catanzaro. Hierarchical Multi-Scale Attention for Semantic Segmentation, 2020. [accessed on 2021-01-26].
- [6] Christopher De Sa, Matthew Feldman, Christopher Ré, and Kunle Olukotun. Understanding and Optimizing Asynchronous Low-Precision Stochastic Gradient Descent. In *Proceedings of the 2017 ACM/IEEE 44th Annual International Symposium on Computer Architecture (ISCA)*. ACM, 2017.
- [7] Xiangru Lian, Wei Zhang, Ce Zhang, and Ji Liu. Asynchronous Decentralized Parallel Stochastic Gradient Descent. In *Proceedings of the 35th International Conference on Machine Learning (ICML)*, pages 3043–3052. PMLR, 2018.
- [8] Shanshan Zhang, Ce Zhang, Zhao You, et al. Asynchronous Stochastic Gradient Descent for DNN Training. In *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 6660–6663, 2013.
- [9] Wei Zhang, Suyog Gupta, Xiangru Lian, and Ji Liu. Staleness-aware Async-SGD for Distributed Deep Learning, 2016. [accessed on 2021-01-26].

- [10] Nikolay Bogoychev, Marcin Junczys-Dowmunt, Kenneth Heafield, and Alham Fikri Aji. Accelerating Asynchronous Stochastic Gradient Descent for Neural Machine Translation, 2018. [accessed on 2021-01-26].
- [11] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, et al. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- [12] Martín Abadi, Ashish Agarwal, Paul Barham, et al. TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems, 2015. Software available from tensorflow.org, [accessed at 2020-02-24].
- [13] Message Passing Interface Forum. *MPI: A Message-Passing Interface Standard, Version 3.1*. High Performance Computing Center Stuttgart (HLRS), 2015.
- [14] Yuichiro Ueno and Rio Yokota. Exhaustive Study of Hierarchical AllReduce Patterns for Large Messages Between GPUs. In *2019 19th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGRID)*, pages 430–439. IEEE, 2019.
- [15] Hiroaki Mikami, Hisahiro Suganuma, Pongsakorn U.-Chupala, et al. Massively Distributed SGD: ImageNet/ResNet-50 Training in a Flash, 2018. [accessed on 2021-01-26].
- [16] Alexander Sergeev and Mike Del Balso. Horovod: fast and easy distributed deep learning in TensorFlow, 2018. [accessed on 2021-01-26].
- [17] Aaron Clauset. A brief primer on probability distributions, 2011. [accessed on 2021-01-26].
- [18] Ang Li, Shuaiwen Leon Song, Jieyang Chen, et al. Evaluating Modern GPU Interconnect: PCIe, NVLink, NV-SLI, NVSwitch and GPUDirect. *IEEE Transactions on Parallel and Distributed Systems (TPDS)*, 31(1):94–110, Jan 2020.
- [19] Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. QSGD: Communication-Efficient SGD via Gradient Quantization and Encoding, 2017. [accessed on 2021-01-26].
- [20] Markus Götz, Charlotte Debus, Daniel Coquelin, et al. HeAT – a Distributed and GPU-accelerated Tensor Framework for Data Analytics. In *Proceedings of the 19th IEEE International Conference on Big Data (BigData)*, pages 276–288. IEEE, December 2020.
- [21] Shen Li, Yanli Zhao, Rohan Varma, et al. PyTorch Distributed: Experiences on Accelerating Data Parallel Training, 2020. [accessed on 2021-01-26].
- [22] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, et al. ImageNet: A Large-scale Hierarchical Image Database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255. IEEE, June 2009.
- [23] Yuxin Wu, Alexander Kirillov, Francisco Massa, et al. Detectron2. <https://github.com/facebookresearch/detectron2>, 2019.
- [24] Marius Cordts, Mohamed Omran, Sebastian Ramos, et al. The Cityscapes Dataset for Semantic Urban Scene Understanding. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3213–3223. IEEE, June 2016.
- [25] Dorian Krause. JUWELS: Modular Tier-0/1 Supercomputer at the Jülich Supercomputing Centre. *Journal of Large-scale Research Facilities*, 5:135, 2019.
- [26] NVIDIA Corporation. NVIDIA Data Loading Library (DALI), 2021. [accessed on 2021-01-25].
- [27] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, et al. Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour, 2018. [accessed on 2021-01-26].
- [28] Hamid Rezatofighi, Nathan Tsoi, JunYoung Gwak, et al. Generalized Intersection over Union: A Metric and A Loss for Bounding Box Regression, 2019. [accessed on 2021-01-26].
- [29] Shuai Zhao, Yang Wang, Zheng Yang, and Deng Cai. Region Mutual Information Loss for Semantic Segmentation, 2019. [accessed on 2021-01-26].
- [30] Andy B Yoo, Morris A Jette, and Mark Grondona. SLURM: Simple Linux Utility for Resource Management. In *Workshop on Job Scheduling Strategies for Parallel Processing*, pages 44–60. Springer, 2003.
- [31] L. Bottou, Frank E. Curtis, and J. Nocedal. Optimization Methods for Large-Scale Machine Learning. *ArXiv*, 2018. [accessed on 2021-01-26].

A Supplementary Materials

A.1 Proof of Convergence

Proof. The following proof of DASO’s global synchronization method is based heavily on the convergence analysis shown by [31] and will show that the gradients determined with DASO are bounded.

Let $X \subset \mathbb{R}^n$ be a known set, and $f : X \rightarrow \mathbb{R}$ a differentiable, convex, L -smooth, and unknown function. Then, the estimator of the stochastic gradient of $f(x)$ is a function $\tilde{g}(x)$ for inputs x determined by the realization of a random variable ζ , such that $\mathbb{E}[\tilde{g}(x; \zeta)] = \nabla f(x : \zeta)$. In the following, ζ is omitted due to space constraints. The stochastic gradient descent (SGD) algorithm updates a model’s state at batch $t + 1$, x_{t+1} , with the following rule $x_{t+1} = x_t - \eta \tilde{g}(x_t)$, where η is the parametric learning rate.

A commonly used variant of SGD in practice is minibatching for computational efficiency reasons. In minibatch SGD, the true stochastic gradient is approximated by averaging across m input items x_i , i.e. $\tilde{G}(x_t) = \frac{1}{m} \sum_{i=1}^m \tilde{g}(x_{t,i})$. The model state x_{t+1} for minibatch SGD is

$$x_{t+1} = x_t - \eta \tilde{G}(x_t) \quad (2)$$

where $\tilde{G}(x_t)$ is an estimator of $\nabla f(x_t)$.

Let us now consider, that S subsequent update steps are performed. It is possible to write the model state as:

$$x_{t+S} = x_t - \eta \sum_{i=0}^{S-1} \tilde{G}(x_{t+i}) \quad (3)$$

One of the primary assumptions in SGD is the Lipschitz-continuous objective gradients. This has the effect that:

$$\begin{aligned} f(x_{t+1}) - f(x_t) &\leq -\eta \nabla f(x_t)^T \mathbb{E}[\tilde{g}(x_t)] \\ &\quad + \frac{1}{2} \eta^2 L \mathbb{E}[\|\tilde{g}(x_t)\|_2^2] \end{aligned} \quad (4)$$

where the Lipschitz constant, L , is greater than zero. Equation (4) implies that the expected decrease in the objective function, $f(x)$, is bounded above by a set quantity, regardless of how the stochastic gradients arrived at x_t [31].

In DASO, the local synchronization step is bound via the same assumptions as minibatch SGD outlined in [31], so long as the iid assumption is upheld. However, the non-standard global synchronization step used in DASO must be shown to be bound under the same principles. DASO’s global synchronization is:

$$x_{t+S}^{\text{DASO}} = \frac{2Sx_{l:t+S-1} + \sum_{i=1}^P x_{p:t}^i}{2S + P} \quad (5)$$

where the l and p subscripts represent the node-local and global model states, S is the number of local update steps before global synchronization, and P is the number of processes.

Similar to Equation (2), this can also be represented via the locally and globally calculated gradients, $\tilde{G}_l(x_{l:t})$ and $\tilde{G}_p(x_{p:t})$ respectively. The global synchronization function in the gradient representation is as follows:

$$x_{t+S}^{\text{DASO}} = x_t - \alpha \left(2S \sum_{k=0}^{S-1} \tilde{G}_l(x_{l:t+k}) + \sum_{i=1}^P \tilde{G}_p(x_{p:t}^i) \right) \quad (6)$$

where $\alpha = \eta/(2S+P)$. Using this, Equation (2), and the fact that the updates between t and S are local synchronizations which take the form of Equation (3), we find that globally calculated gradients are as follows.

$$\begin{aligned} \tilde{G}^{\text{DASO}}(x_{t+S-1}) &= P \sum_{\beta=0}^{S-1} \tilde{G}_l(x_{l:t+S-\beta}) \\ &\quad - 2S \tilde{G}_l(x_{l:t+S-1}) \\ &\quad + \sum_{i=1}^P \tilde{G}_p(x_{p:t}^i) \end{aligned} \quad (7)$$

As all gradient elements in Equation (7) are bound under Equation (4), $\tilde{G}^{\text{DASO}}(x_{t+S-1})$ is similarly bounded. \square