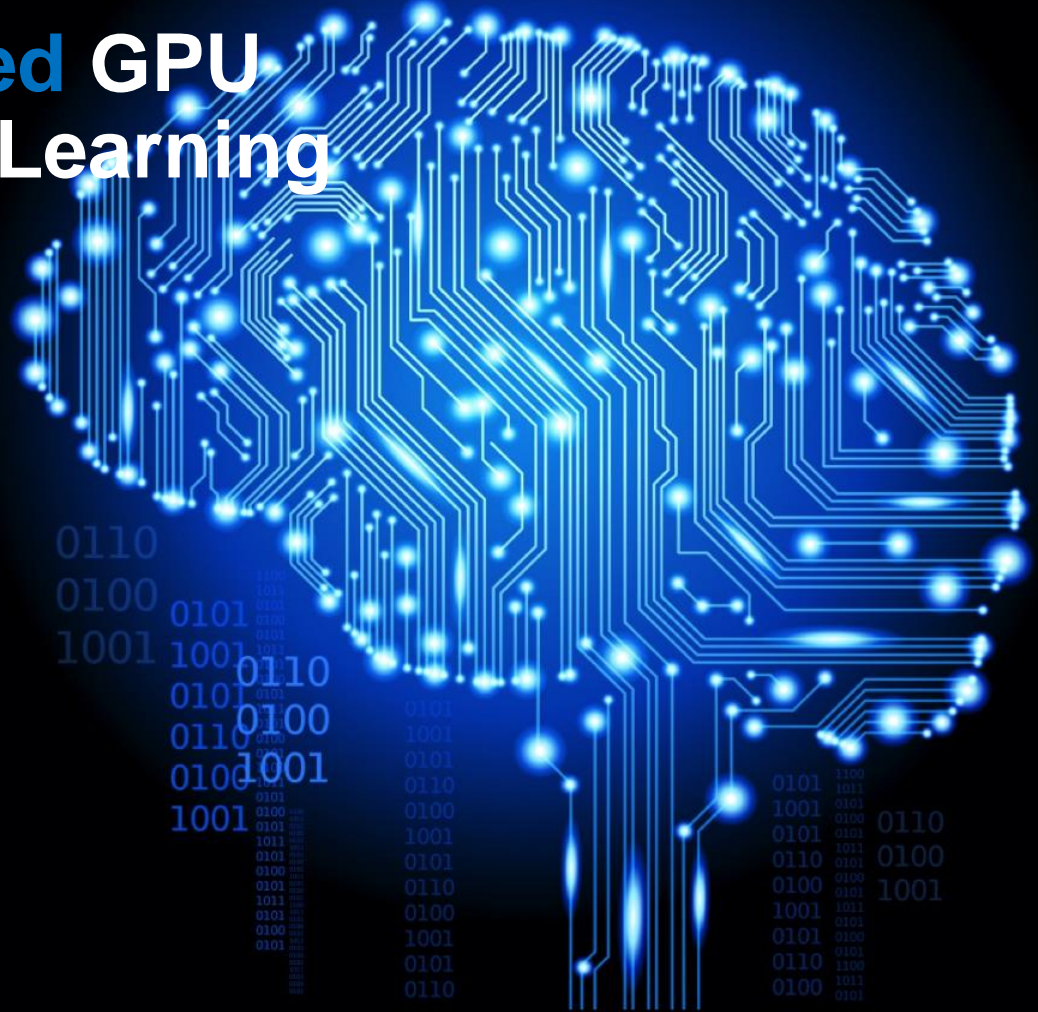


Optimizing distributed GPU collectives for Deep Learning Workloads



Pidad D'Souza (pidsouza@in.ibm.com)

Nysal Jan K A (inysal@in.ibm.com)

ISDL, IBM India Pvt. Ltd

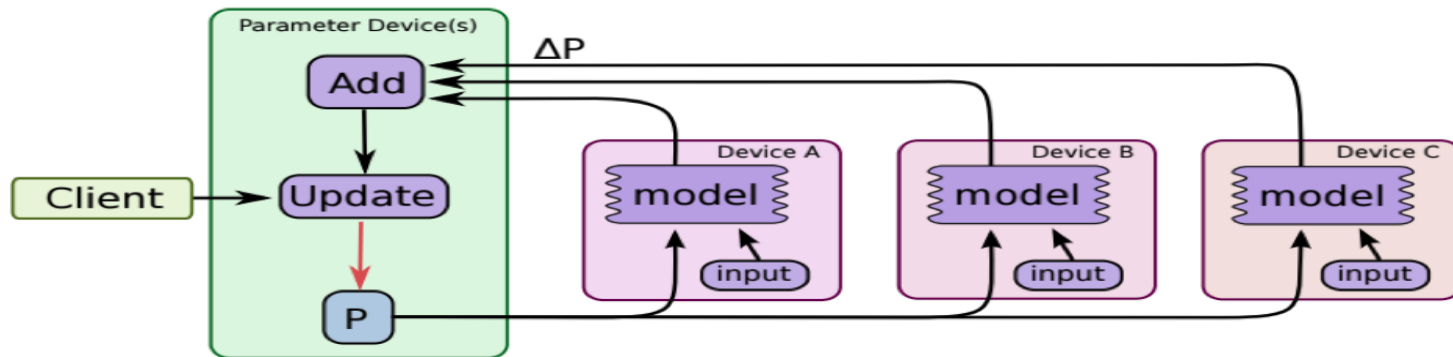
Agenda

- Distributed Training
- Neural network in Deep Learning
- Baidu Allreduce
- GPU Collective
- Observations
- Future work

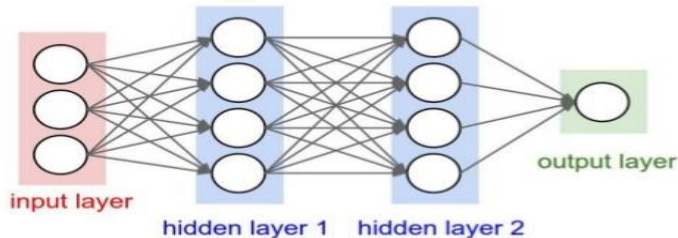
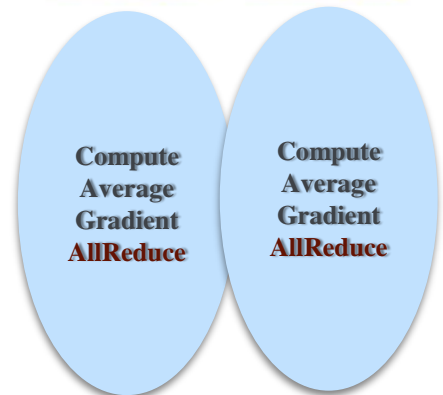
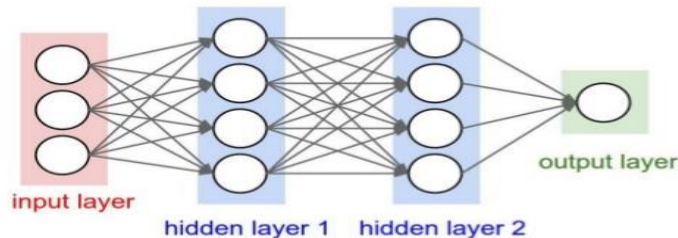
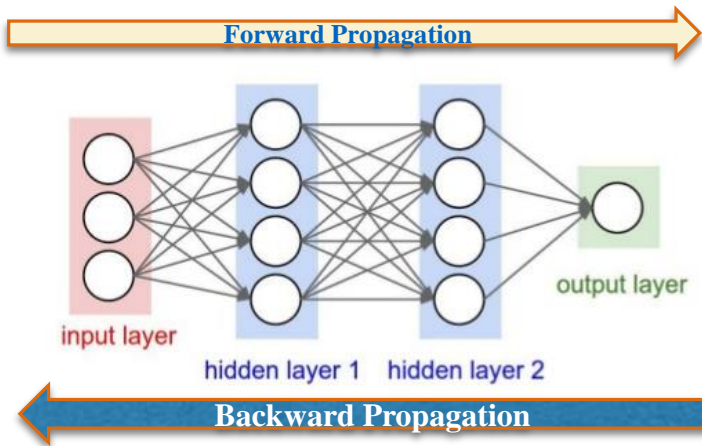


Distributed Training

- Deep learning: Image/face recognition, object classification, language processing, sign board identification etc.
 - Compute intensive training stage – accelerate(GPUs)
 - Ever increasing size of datasets and large number of parameter sets
 - Run models on multiple GPUs to speed up Training stage
 - Achieve higher accuracy of prediction
- Framework's like Caffe2, Microsoft's CNTK and Amazon's DSSTNE
 - Exploit multiple GPUs across multiple nodes using existing MPI runtimes
 - Message sizes ~256MB and more (varies with the model)



Neural Networks in Deep Learning



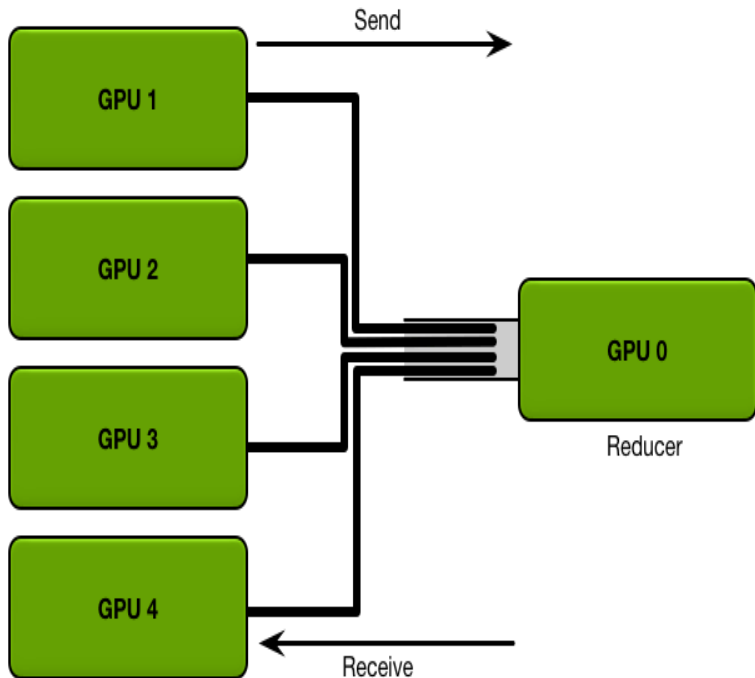
Need for Distributed ML

- Large number of parameters with multiple Neural network layers

Each Training Iteration

- Forward Propagation
- Backward Propagation
- Update Weights

Prior-art: Baidu Ring Allreduce



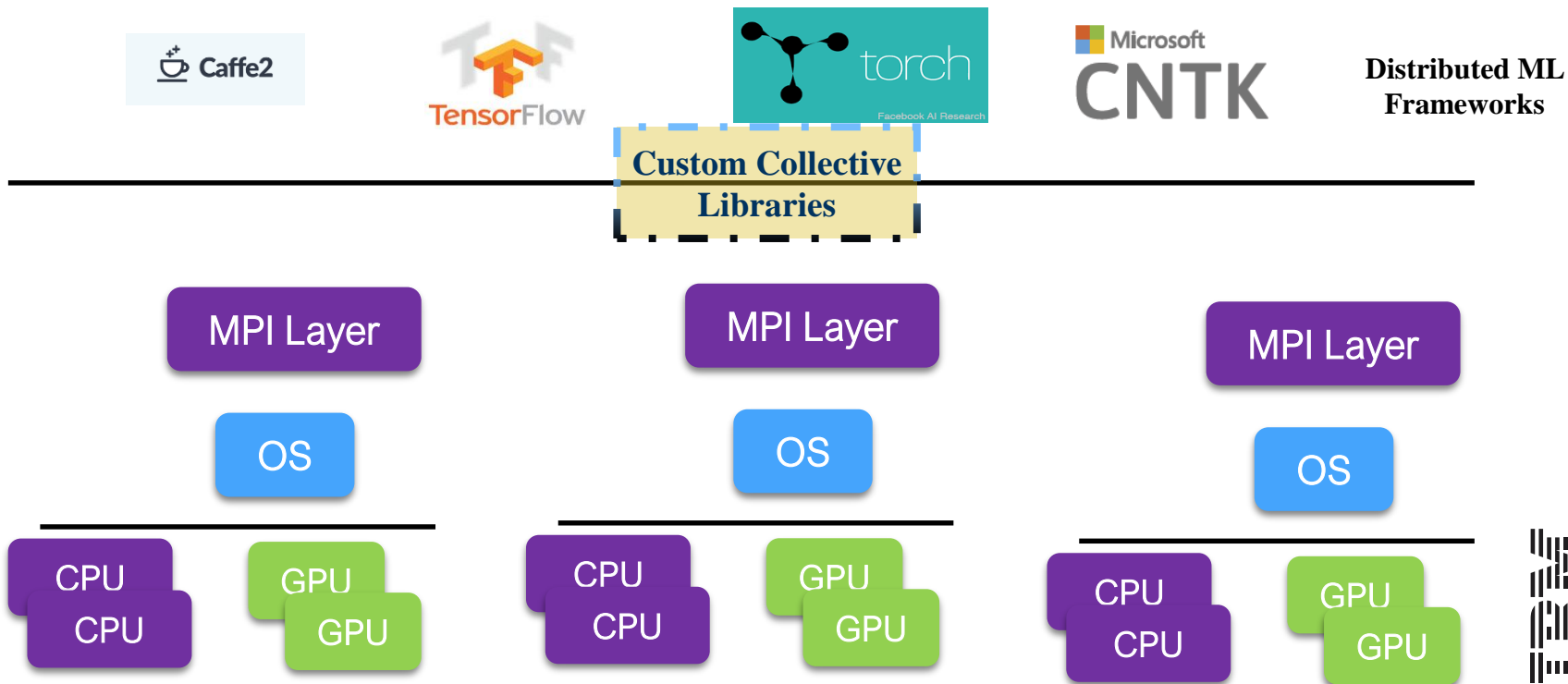
Data transfer to and from a single reducer GPU

- Distribute different operations onto GPUs through data parallel stochastic gradient descent (SGD)
- Subset of the samples in the minibatch are assigned to each GPU
- GPUs communicate with each other to average the gradients
- Apply the averaged gradient to the weights to obtain new weights

Ref: <http://research.baidu.com/bringing-hpc-techniques-deep-learning/>



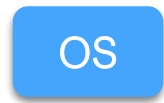
Distributed Machine Learning Stack



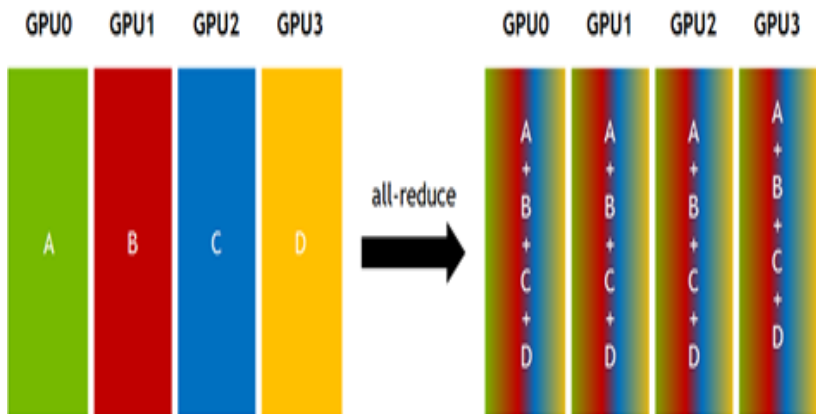
GPU Collective component in Distributed Machine Learning Frameworks



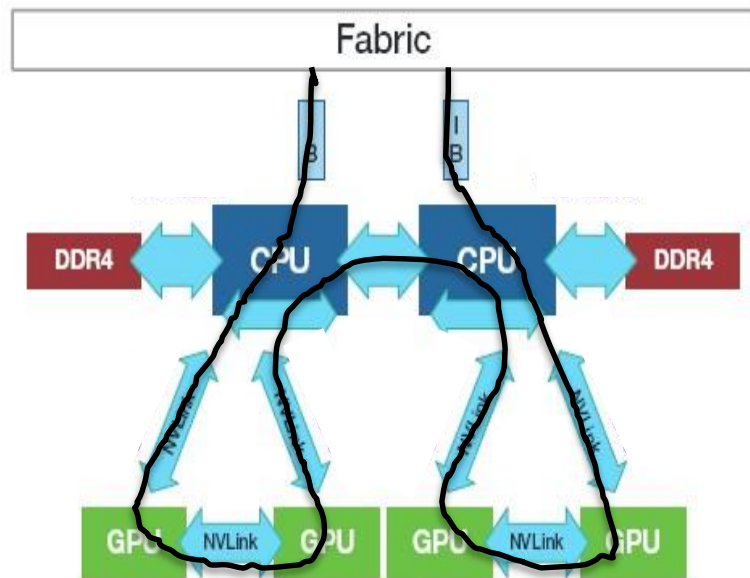
Distributed ML Frameworks



All Reduce

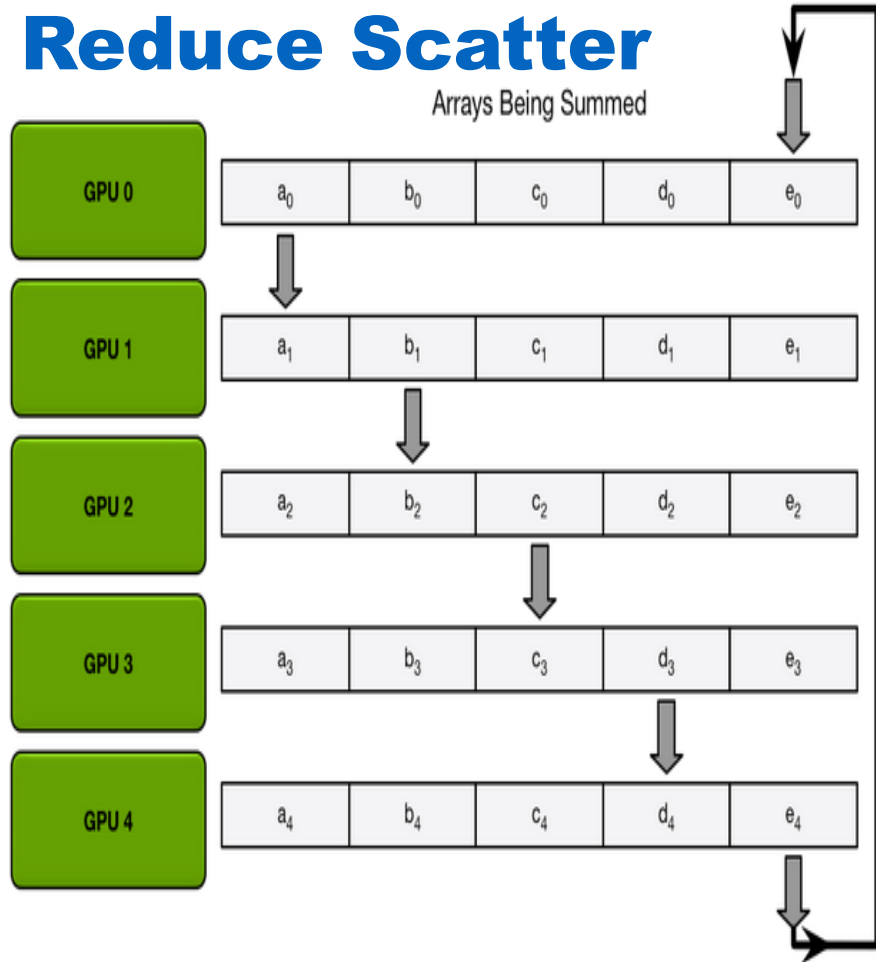


Ring Topology

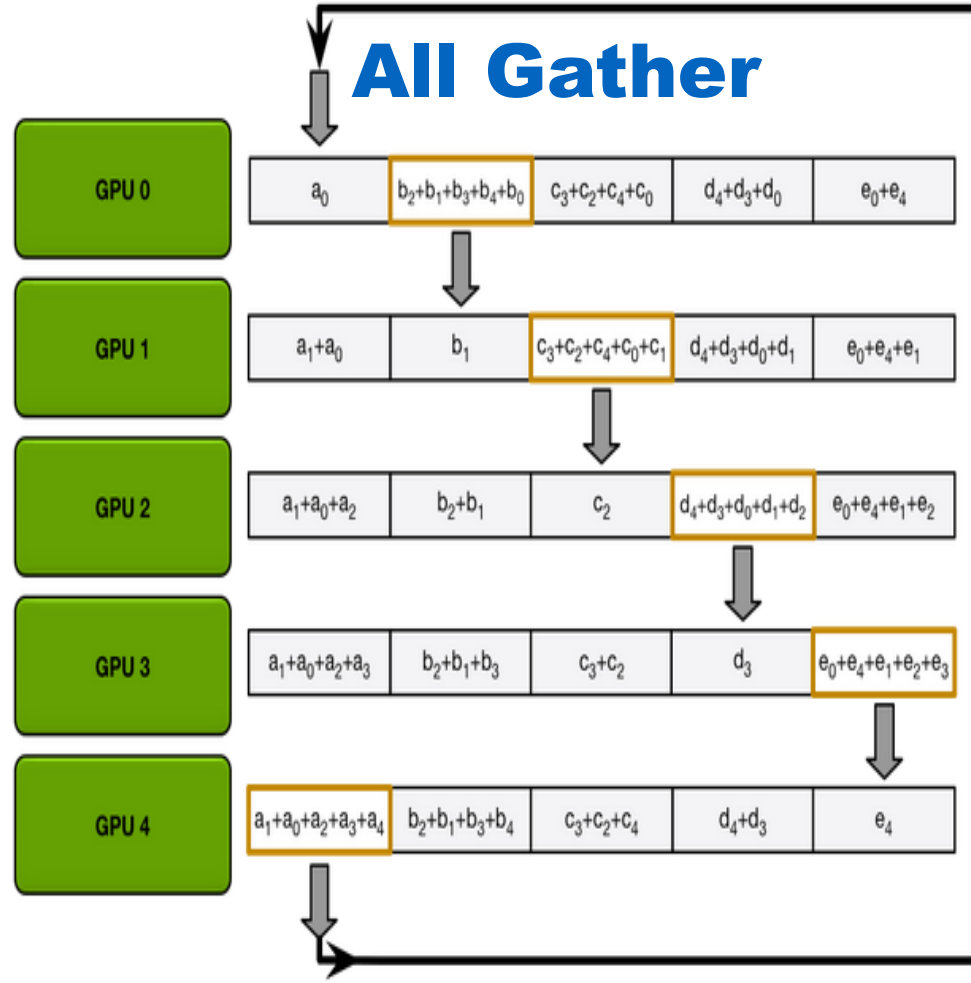


Reduce Scatter

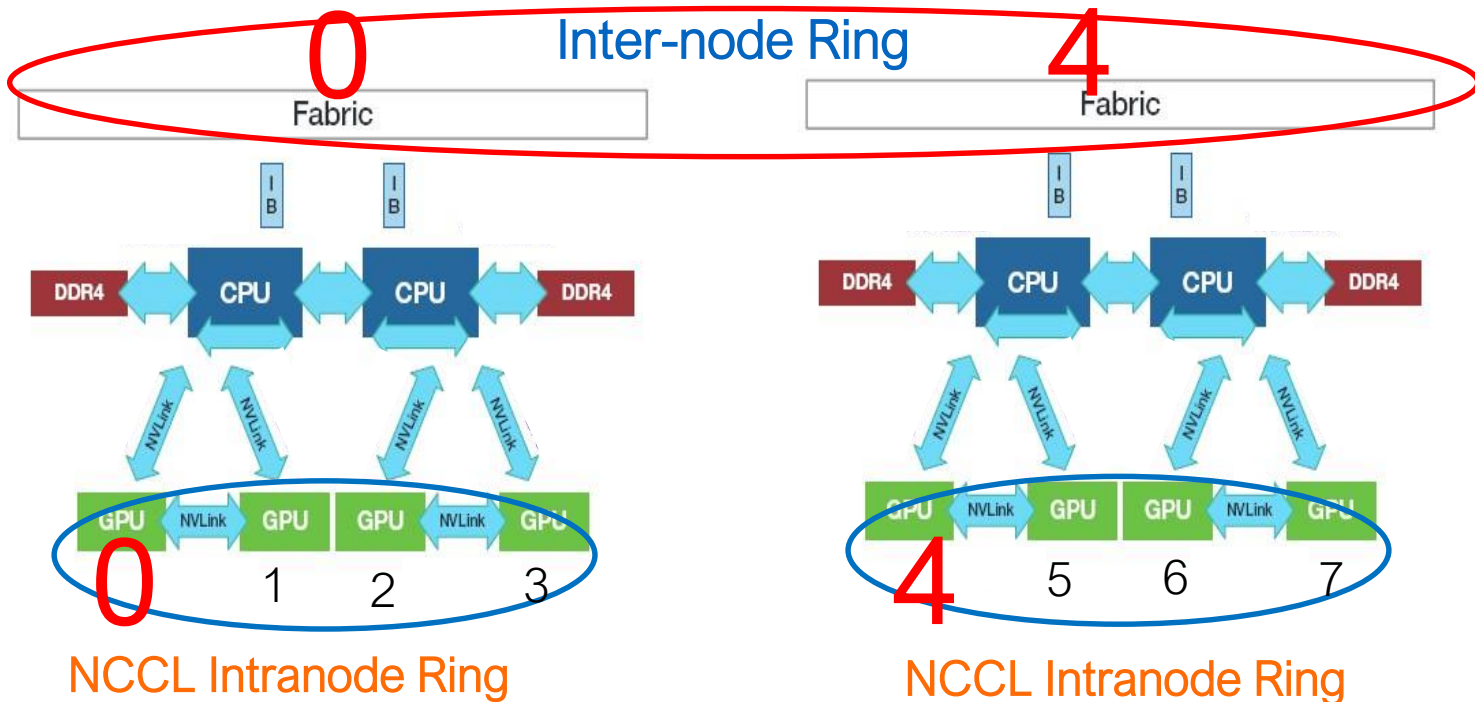
Arrays Being Summed



All Gather



Leader based hierarchical Allreduce Algorithm



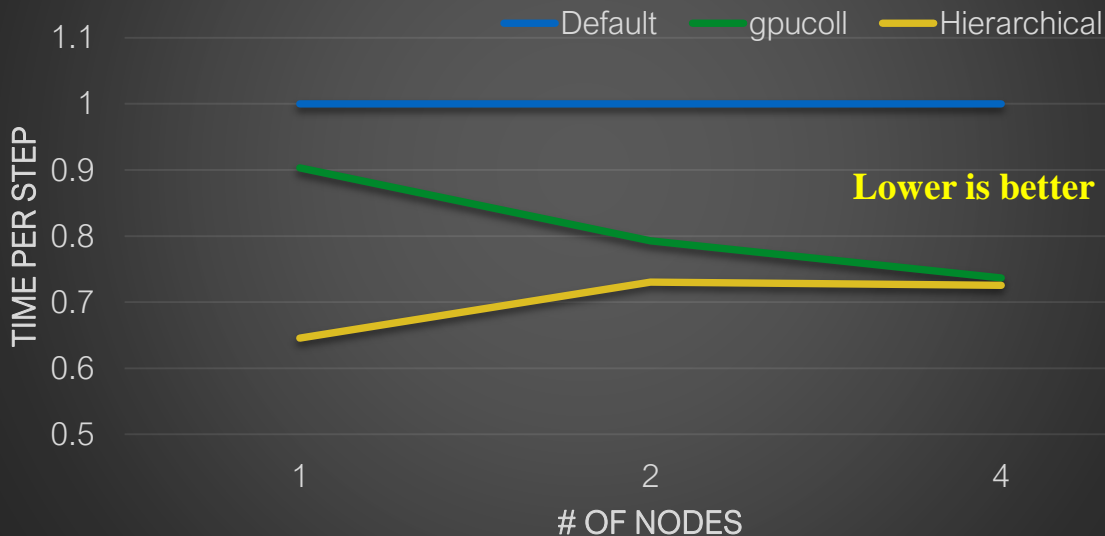
NCCL used for demonstration purpose only, GPU Collective component can have it's own implementation for intra-node collective.



TensorFlow Inception-v3

with ImageNet Dataset (on IBM POWER S822LC for HPC)

Normalized Time per step Scaling on multiple nodes



GPU Coll vs Default Gain

1 Node: 9.7%

2 Nodes: 20.7%

4 Nodes: 26.4%

GPU H-Coll vs Default Gain

1 Node: 35.4%

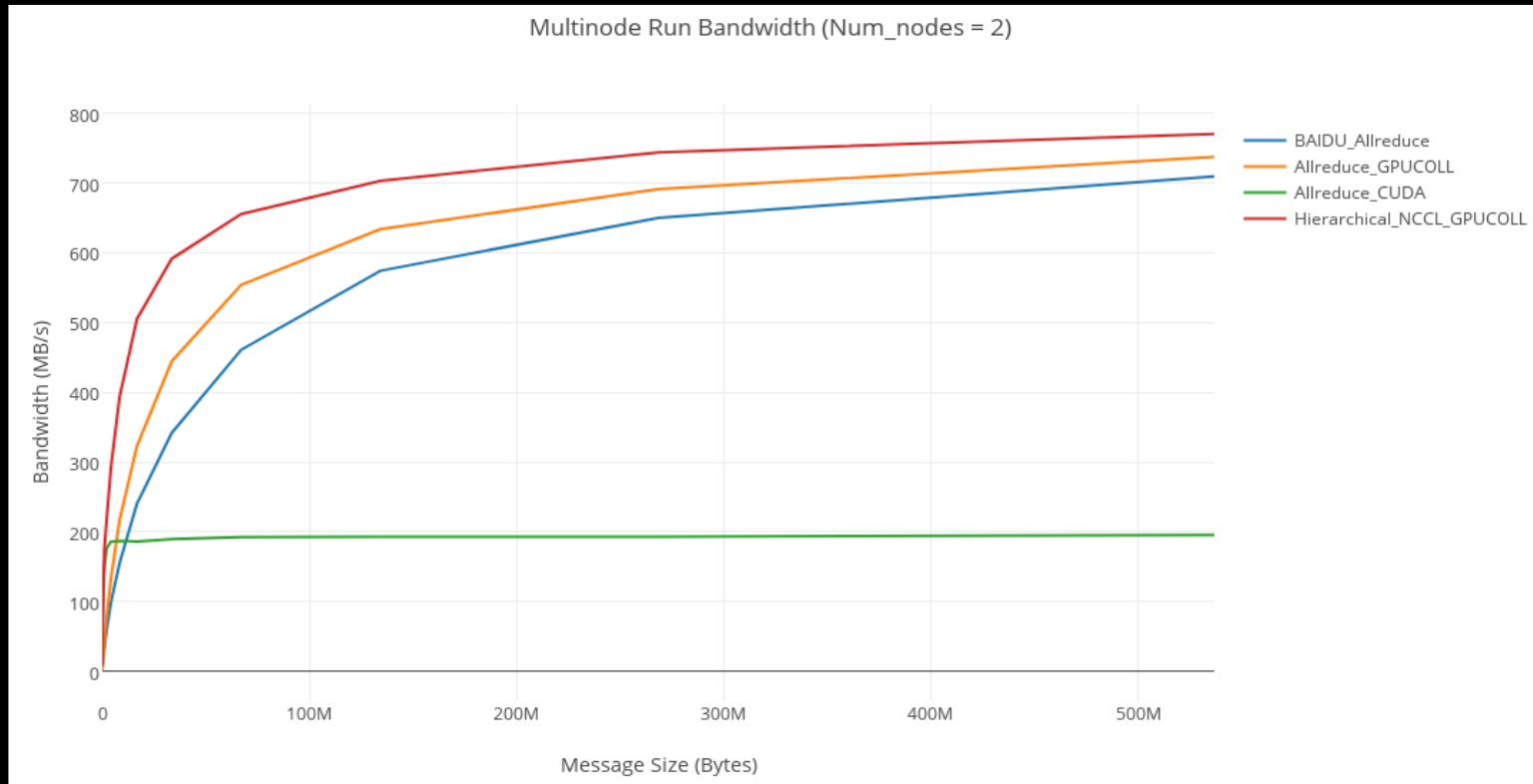
2 Nodes: 27%

4 Nodes: 27.5%

Additional optimization should help improve hierarchical implementation performance.



Bandwidth Benchmark Results



Note : Proof points demonstrated in this chart are for reference only, to demonstrate advantages of newer implementation and not the best achievable results.



Future Work

- Multi-Leader Approach
- Explore GPU Direct RDMA on IBM POWER9
- GPU Async Support
- Network Topology Based Optimization



Conclusion

- Performance improvements with GPU Collective implementations
- Frameworks can seamlessly take advantage, without integration efforts

For a comparison of NCCL, GLOO, and all-reduce performance to PowerAI DDL's leadership performance, see Minsik Cho, "Efficient Communication Library for Large-Scale Deep Learning," (S8479) GTC 2018

Acknowledgement

Mayank Roy, IIT – Kharagpur,
Abhishek Tiwari, IIT – Kharagpur



Thank You

