



Hy-Fi: <u>Hy</u>brid <u>Fi</u>ve-Dimensional Parallel DNN Training on High-Performance GPU Clusters

Presentation at MUG '22

Arpan Jain

Network Based Computing Laboratory (NBCL)

Dept. of Computer Science and Engineering, The Ohio State University

jain.575@osu.edu

- Introduction and Motivation
- Problems and Challenges
- Key Contributions
- Performance Evaluation
- Conclusion

Deep Learning meets Super Computers

 NVIDIA GPUs - major force for accelerating DL workloads



ISC '22

Performance Share 7.5% 7.8% 59.5% 8.9% NVIDIA Tesla V100 NVIDIA A100 NVIDIA Tesla V100 SXM2 NVIDIA A100 80GB NVIDIA A100 40GB NVIDIA A100 SXM4 40 GB NVIDIA Tesla P100 NVIDIA Volta GV100 NVIDIA Tesla K40 Matrix-2000 Others www.top500.org

Accelerator/CP Family

Courtesy: https://openai.com/blog/ai-and-compute/ Network Based Computing Laboratory

High-Performance Deep Learning

Parallelization Strategies

- Data Parallelism
- Model Parallelism
 - **Distributed Training** Layer-level Parallelism Layer Data Parallelism Model Parallelism Hybrid Parallelism ٠ Pipeline Sub-Graph ٠ Advance Offload **Neuron-level Parallelism** Schemes Layer-level Parallelism **Neuron-level Parallelism** D&SP Megatron Spatial Channel ٠ Hybrid Parallelism Sub-Graph Layer Pipeline Spatial Channels Hy-Fi Parallelism Parallelism Parallelism Parallelism Parallelism D&SP
 - Megatron
 - Hy-Fi

٠

Why Hybrid Parallelism?

- Data-Parallelism only for models that fit the memory
- Out-of-core models
 - Deeper model
 Better accuracy but more memory required!
 - Real-world applications have very highresolution images
- Model parallelism can work for out-ofcore models!
 - Performance is questionable!
 - Layer-parallelism is not enough



Square Image Size (X * X)

- Introduction and Motivation
- Problems and Challenges
- Key Contributions
- Performance Evaluation
- Conclusion

Research Challenges

Challenge-1: Halo Exchange in PyTorch Challenge-2: Exploitation of different parallelism dimension

Challenge-3: Scaling Integrated Hybrid Training Solutions

Meet Hy-Fi!

- Introduction and Motivation
- Problems and Challenges
- Key Contributions
- Performance Evaluation
- Conclusion

Proposed Hybrid Five-Dimensional Parallelism (Hy-Fi)

• Integrates spatial, layer, pipeline, bi-directional, and data parallelism



Strategies to Optimize Halo Exchange

- Halo Exchange
 - Needed to compute convolution and pooling operations
- Proposed Halo-D2 and explored different image distribution strategies





X Convolution Halo Exchange o Data locally available Φ

#: Process Number 1 2 3

7

11

15

6

10

14

5

9

13

Image Distribution Strategies #: Process Number

Conv

2

Conv

2

Conv

2

Conv

2

Conv

3

Conv

3

Conv

3

Conv

3





Conv

Λ

Conv

4

Conv

4

Conv

4

Halo

Ex.

n=3

Conv

Conv

Conv

5

Conv

5

Conv

Conv

Conv



4

8

12

16

Halo Ex. n: Number of columns/rows to exchange

Communication using send/recv operations

Proposed Communication OptimizationHalo-D2

Network Based Computing Laboratory

ISC '22

High-Performance Deep Learning

GEMS-MAST vs GEMS-MASTER

- Bi-directional Parallelism
 - Hard to overlap Allreduce operation ->
- Proposed communication optimization





- Introduction and Motivation
- Problems and Challenges
- Key Contributions
- Performance Evaluation
- Conclusion

Evaluation Setup

- System
 - Lassen at Lawrence Livermore National Laboratory (LLNL)
 - POWER9 processor
 - 4 NVIDIA Volta V100 GPUs per node
- Interconnect
 - X Bus to connect two NUMA Nodes
 - NVLink is used to connect GPU-GPU and GPU-Processor
 - Infiband EDR
- PyTorch, MVAPICH2-GDR 2.3.5
- We use and modify model definitions for ResNet(s) from *keras.applications* and AmoebaNet model from *TorchGpipe*

Accelerating Out-of-core Training at Scale

- Approach
 - LP: Layer Parallelism
 - Pipeline: Pipeline Parallelism
 - SP: Spatial Parallelism
 - Master: Hy-Fi
 - Opt: with Communication optimization in Hy-Fi
- Setup
 - Image Size: 2048 X 2048
 - 2048 NVIDIA V100 GPUs
- Speedup
 - Up to 2.67X over Layer Parallelism (LP)
 - Near-linear scaling (94.5%) on 2,048 GPUs



AmobeNet-f416

Enabling Training on Very High-Resolution Images

- Enabled training on 8,192 X 8,192 and 16,384 X 16,384 images sizes
- Speedup over basic spatial parallelism
 - 8,192 X 8,192 Images: 1.476X and 2.26X (Strong Scaling)
 - 16,384 X 16,384 Images: 1.47X



AmobeNet-f214 on 16,384 X 16,384 images



AmobeNet-f214 on 8,192 X 8,192 images

ISC '22

- Introduction and Motivation
- Problems and Challenges
- Key Contributions
- Performance Evaluation
- Conclusion

Conclusion

- Proposed and Designed Hy-Fi
 - Integrated five different parallelization strategies
 - Spatial, Layer, Pipeline, Bi-directional, and Data
 - Communication optimizations to improve speedup
 - PyTorch and MPI for flexibility and scalability
- Performance Evaluation on large systems
 - Up to 2.67X speedup for out-of-core DNNS
 - Scaled Hy-Fi to 2,048 V100 GPUs on LLNL Lassen
 - Achieved 94.5% scaling efficiency with GEMS-Hybrid
- Future Work
 - Use Hy-Fi to train out-of-core DNNs on larger image tiles for digital pathology
 - Add more parallelization strategies

Thank You!

jain.575@osu.edu

Network-Based Computing Laboratory http://nowlab.cse.ohio-state.edu/

High Performance Deep Learning <u>http://hidl.cse.ohio-state.edu/</u>



The High-Performance Deep Learning Project <u>http://hidl.cse.ohio-state.edu/</u>



The High-Performance MPI/PGAS Project http://mvapich.cse.ohio-state.edu/

Network Based Computing Laboratory



High-Performance Deep Learning