



#### Accelerating DNN training on BlueField DPUs

#### Presentation text MUG '21

Arpan Jain

Network Based Computing Laboratory (NBCL) Dept. of Computer Science and Engineering , The Ohio State University Jain.575@osu.edu

## **BlueField DPU / Smart NIC Architecture**

- BlueField includes the ConnectX6 network adapter and data processing cores
- System-on-chip containing 64-bit ARMv8 A72
- Why BlueField DPU for Deep Learning?
  - State-of-the-art DPUs bring more compute power to network
  - Deep Learning training needs all the available compute power it can get





MUG '21

Model Validation Yes **END TRAINING** 

3

**Data Augmentation** 

Training

## Offload Naive (O-N): Offloading DL Training using Data Paralle jan Illelism can used to train DNN on DPUs



A. Jain, N. Alnaasan, A. Shafi, H. Subramoni, D. Panda, "Accelerating CPU-based Distributed DNN Training on Modern HPC Clusters using BlueField-2

Laboratory<sup>128</sup>

#### MUG <u>'21</u>

#### High-Performance Deep Learning

## **Accelerating DNN Training using Offload-**

**Naive**per iteration can be used to distribute the work (batch size) between CPU and DPU

- Speedup:
  - We report up to 1.03X speedup
  - Maximum speedup possible: 1.04X
- Offload-Naive does not give significant speedup as forward and backward pass are computeintensive tasks and DPUs are not as powerful as CPUs



A. Jain, N. Alnaasan, A. Shafi, H. Subramoni, D. Panda, "Accelerating CPU-based Distributed DNN Training on Modern HPC Clusters using BlueField-2

## **Design 1: Offload Data Augmentation (O-DA)**

**MUG '21** 

- Offloads the reading of training data from memory and data augmentation on input data to DPUs.
- Creates two types of processes
  - Training processes (on CPU)
  - Data Augmentation processes (On DPU)
- Initializes two buffers to enable asynchronous communication
- Each training processes has one data augmentation processes on

Laboratory

A. Jain Di Palueasan, A. Shafi, H. Subramoni, D. Panda, "Accelerating CPU-based Distributed DNN Training on Modern HPC Clusters using BlueField-2 DPUs",



High-Performance Deep Learning



#### **Design 2: Offload Model Validation (O-MV)**

## **Design 3: Offload Hybrid (O-Hy)**

- Offloads data augmentation and model validation to DPUs.
- Creates three types of processes
  - Training processes (on CPU)
  - Augmentation processes Data (On DPU)
  - Testing processes (On DPU)
- Each Data Augmentation process on supports DPU multiple training processes.
- Data Augmentation processes does asynchronous communication and



### **Training ResNet-20 on CIFAR-10 Dataset**

- Speedup
  - Single node: O-DA (13.8%) and O-MV (3.1%)
  - Multi-node: Achieves average 13.9% speedup on 1-16 nodes



Laboratory<sup>28</sup>

MUG '21

### **Training ResNet-56 on SVHN Dataset**

• Speedup

Single node: O-DA (7%), O-MV (5.5%), and O-Hy (10.1%)



A. Jain, N. Alnaasan X peri mentenii, D. Panda, "Accelerating CPU-based Distributed DNN Training Operimentes using BlueField-2

Laboratory<sup>28</sup>

MUG '21

High-Performance Deep Learning

## **Training ShuffleNet on Tiny ImageNet Dataset**

- Speedup
  - Single node: O-DA (12.5%), O-MV (1.2%), and O-Hy (8.9%)
  - Multi-node: 10.2% speedup on 16 nodes



Laboratory<sup>28</sup>

**MUG '21** 

High-Performance Deep Learning 11

### Conclusion

- Proposed novel offloading designs for DPUs
  - Offload Naive
  - Offload Data Augmentation
  - Offload Model Validation
  - Offload Hybrid
- Reported up to 15%, 12.5%, and 11.2% speedup for CIFAR-10, SVHN, and Tiny ImageNet datasets
- Demonstrated consistent performance gain on multiple nodes.
- Uses Torchvision, PyTorch, Horovod, and MPI for flexibility and scalability
- Future Work
  - Use DPUs to accelerate DNN training on GPUs
  - Evaluate TransFormer models

# **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>

Laboratory



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

High-Performance Deep Learning 13

