

National Institute of Advanced Industrial Science and Technology (AIST), Japan

# Large-Scale Design of Distributed Deep Learning

#### **Truong Thao Nguyen**



Digital Architecture Research Center Continuum Computing Architecture Research Team



AIST-Tokyo Tech Real World Big-Data Computation Open Innovation Laboratory (RWBC-OIL)

29<sup>th</sup> July 2021

# Introduction

### **Self Introduction**

#### **Personal information**

- TRUONG THAO NGUYEN
- Email: <u>nguyen.truong@aist.go.jp</u>, <u>nguyen.88.tt@gmail.com</u>
- Web: <u>https://researchmap.jp/nguyentt?lang=en</u>

#### **Research expertise:**

Interconnection network, HPC System, Distributed computing, High-Performance Large-scale Distributed Deep Learning



- Design Interconnection Network for High Performance Computing System.
  - Master at SoICT/HUST, VietNam (2012-2014)
  - PhD at SOKENDAI / NII, Japan (Apr 2015 Mar 2018)
- High-Performance Large-scale Distributed Deep Learning.
  - Postdoc at RWBC-OIL, AIST (Apr 2018 Sept 2019)
  - Researcher at AIST (Oct 2019 present)



- <u>https://www.aist.go.jp/</u>
- The National Institute of Advanced Industrial Science and Technology (AIST), one of the largest public research organizations in Japan, focuses on the creation and practical realization of technologies useful to Japanese industry and society, and on **"bridging" the gap between innovative technological seeds and commercialization**.
- AIST is organized into **5 departments and 2 centers** that bring together core technologies to exert its comprehensive strength.
- AIST is also actively building a global network by, for example, signing memorandums of understanding for comprehensive research cooperation (MOUs) with major research institutes around the world.

#### **Information Technology and Human Factors**

From IT Infrastructure to application technology: make, connect, use

Research Unit

- Human Informatics and Interaction Research Institute
- Cyber Physical Security Research Center
- Human Augmentation Research Center
- Human-Centered Mobility Research Center
- Artificial Intelligence Research Center
- Industrial Cyber-Physical Systems Research Center
- Digital Architecture Research Center

Open Innovation Laboratory

https://www.aist.go.jp/aist\_e/dept/en\_dithf.html The rapid growth of devices (smartphone, tablets, robots, wearables, etc.) and the Internet has increased the amount of information that is being produced and accessed by society. In order to better utilize the data produced from millions of devices and systems, we are conducting research and development in a wide range of fields at the interface between information technologies and human factors. Our mission is to engage and enrich the public through the research and development of intelligent systems combining computational and physical capabilities for human use.

# **ABCI - The World #12 Supercomputer**



AIRC Artificial Intelligence Cloud Research Team(~ 2021/03/31) <u>https://www.airc.aist.go.jp/en/teams/</u>

AI Bridging Cloud Infrastructure (ABCI) is the world's first large-scale Open AI Computing Infrastructure

ABCI consists of; 120 Compute Nodes (A) that form in total 960 NVIDIA A100 GPU accelerators, 1,088 Compute Nodes (V) that form in total 4,352 NVIDIA GPU V100 accelerators, shared file systems that provide in total 35PB capacity, Cloud Storage,







#### 👑 Digital Architecture Research Center

The Digital Architecture Research Center was established to develop and implement technologies that promote the digital revolution as a contribution to solving social issues and innovative innovation. New services and value for all individuals and organizations by appropriately linking various, diverse, and large amounts of data generated by individuals, devices, factories, etc., and highly integrating cyberspace and physical space while utilizing AI technology. We aim to realize an advanced society that continuously provides new services. Coordinating research units inside and outside the domain to advance the data-driven digital society, it acts as a control tower for architectural design, necessary technology development, standardization, software implementation and dissemination activities.

#### https://www.digiarc.aist.go.jp/en/org/

Continuum Computing Architecture Research Team

Leader Yusuke Tanimura



Continuum Computing Infrastructure Research Team

Leader Takahiro Hirofuchi

# This Talk - Overview

### **Motivations**

- Explosion of Deep Learning (DL)
  - Effectiveness in a variety of applications
- Long training time limits the development of new DL applications
  - Training GPT-3 model takes <u>355 years</u> on a V100 GPU server and <u>cost \$4.6M</u>
- Parallel training on Supercomputer / HPC systems
- $\rightarrow$  Our target: large-scale training, e.g., <u>100s-1000s of GPUs</u>



Illustration of data parallelism with Supercomputer

Gaming

Music

Technology

Medical Science

**Computer Vision** 

**Robotics** 

Communication

Programming

# Long Training Time Limits the Development of New Deep Learning Applications



### **Challenges of Training DNN in New Domains**



### **HOW - Research Goal and Approach**

Goal: enable training DNNs on large-scale HPC system in the magnitude of hours.

Technical Goal: 1 Less Memory 2 Faster 3 Accuracy



# High Performance I/O

### **HOW - Research Goal and Approach**

Goal: enable training DNNs on large-scale HPC system in the magnitude of hours.

Technical Goal: 1 Less Memory 2 Faster 3 Accuracy



16

# **Bigger Dataset - I/O Challenges**

- Random access with Global Shuffling each epoch
- Number of input files is large (millions billions)
  - File size is also big, e.g., Cosmoflow
- Load the whole dataset into local storage before training **Problems**: Dataset size is bigger than local storage (e.g., SSD)
- Bad performance with PFS
- Related work: Local Shuffling







# High Performance Communication

### **Targeted Communication Model**

Synchronous: communication start when the computation are finished at all the computing nodes.



Parameter Server model (Reduce + Broadcast) GPU1 GPU2 ... GPU P-1 GPU P N GPU2 ... GPU P-1 GPU P GPU1 GPU2 GPU2 GPUP-1 GPU P

Distributed model (our target) (Allreduce)

#### **Distributed Communication becomes practical**

### **Limitation of Data Parallelism in Future**





Communication becomes bottleneck

[\*] Dryden, Nikoli, et al. "Aluminum: An asynchronous, GPU-aware communication library optimized for large-scale training of deep neural networks on HPC systems", MLHPC 2018.

#### Large-Scale Distributed Deep Learning



# **Approach 1: Communication Algorithm**

Goal: enable training DNNs on large-scale HPC system in the magnitude of hours.

Technical Goal: 1 Less Memory 2 Faster 3 Accuracy



23

# **Optimizing Message Passing Algorithms**

Problem: High communication time due to (1) large message size and (2) larger scale

Task: Study the method to reduce the communication time while maintain the accuracy



[2] N. T. Truong, M. Wahib, and R. Takano, "Efficient MPI-Allreduce for Large-Scale Deep Learning on GPU-Clusters", in Concurrency and Computation Practice and Experience, Special Issue on Parallel and Distributed Computing and Networking (CCPE) 24

# **Data Compression**

#### **Sub-issues of sparcification:**

(1) <u>Ineffective pair-wise communication algorithm</u>(2) Computation overhead of top-k select algorithm



[3] N.T. Truong, "Topology-aware Sparse Allreduce for Large-scale Deep Learning", 38th IEEE International Performance Computing and Communications Conference (IPCCC), Oct 2019

# **Sparcification with Mixed Data Representation**

#### **Target: Reduce the communicated message size with Data Compression**

Idea: Compress the message

Communication Challenges: Unknown overlapping of the subsets between GPUs

#### Baseline

- Can not use ring-based Reduce
- ➢ Use pair-wise communication



**Cost** Send all data in O(P) steps 2(P-1)  $\alpha$  + 2(P-1)N<sub>2</sub> $\beta$ 

#### **Proposed Method**

• Dynamically switch between sparse data format and dense data format



Fundamental operations: SPLIT, SUM, COPY
 SMALL EXTRA COMPUTATION γ (<10%)</li>

$$\leq 2(P-1) \alpha + \frac{2(P-1)N_1}{P}\beta + \gamma$$



### **Network Congestion**

**Problem**: Number of involved GPUs (P) increase

- Ring-based algorithm  $O(P)\alpha$ : latency factor increase.
- Halving-Doubling algorithms  $O(\log P)\alpha$ : network congestion



Length-2 shortcut

 $1+2^{Q-1}$ 

2<sup>p</sup> downlinks

..... Length-1 shortcut

2<sup>Q</sup> -7

2<sup>Q</sup> -3

2<sup>Q</sup> -2

Length-8 shortcut



Shilfted Halving – Doubling Algorithm



[4] N.T. Truong, et. al. "An Algorithm and System Co-design for Large-Scale Training of Distributed Deep Learning", CCGRID 2021

#### **Proposal:** Distributed Loop Network topology

### **Micro-benchmark Result**

#### Simulation result on ABCI

- ✓ DLN\_Ring slightly better FatTree\_Ring 2-5%
- ✓ DLN\_SHD significantly better than FatTree\_Ring: 31%, 71%, and 90%
- ✓ DLN\_SHD better than FatTree\_HD: ~ 50%
- ✓ DLN\_SHD better than DLN\_HD: ~ 25%



Better

#### **High-bandwidth Link: Hybrid Switching System for DL**

- Small flows use EPS (electric packet switching)
- Big flows use OPS (optical circuit switching) such as data of DL
  - Logical-ring algorithm never change the pattern
  - No-need of reconfiguration



#### Multi Layer Electrical Switch System

**Electrical and Optical Hybrid Switch System** 

[5] <u>N.T. Truong,</u> et. al. "On the Feasibility of Hybrid Electrical/Optical Switch Architecture for Large-Scale Training of Distributed Deep Learning ",PHOTONICS, 2019 31

# High Performance Parallelism Model

### **HOW - Research Goal and Approach**

Goal: enable training DNNs on large-scale HPC system in the magnitude of hours.

Technical Goal: 1 Less Memory 2 Faster 3 Accuracy



### **Memory Really a Limiting Factor?**

- Yes!
- Examples:
  - Larger models:
    - Megatron-LM: 8.3B parameters O(10) GBs
    - ZeRO-2: 170B parameters O(100) GBs
  - Larger Samples: Common in scientific datasets
    - CosmoFlow: ~2GB/sample O(100) GBs for act. 1<sup>st</sup> layer
  - When Increasing mini-batch size
    - Can go up to 127K with methods such as K-FAC

#### **Solutions?**

#### **Use More Memory/ Node or GPUs!**

In-efficient use of resources

-ve



#### **Out-Of-Core Methods**

- More sensitive to co-located jobs
- Tricky when no local DRAM/SSD
- Simplicity of data parallelism
- No changes to model/framework +ve
- No messing with Batch Norm (BN)

#### **Explicit Model-parallelism**

Constructing a cost model that is accurate
High implementation cost -ve

- Analytical: deterministic

+ve

- Low overhead

### **Out-of-core to Overcome Memory Limitation**

- **Problem**: Model size bigger than GPU's memory
  - 100s TB as in GPT-3 vs. 32GB HBM2
- **Proposal**: OOC in DNN is intuitively swapping layer
  - Move finished layer to CPU memory
  - Redundantly recompute in backward pass





[6] Mohamed Wahib, Haoyu Zhang, <u>**Truong Thao Nguyen**</u>, Aleksandr Drozd, Jens Domke, Lingqi Zhang, Ryousei Takano, Satoshi Matsuoka, "Scaling Distributed Deep Learning Workloads beyond the Memory Capacity with KARMA", Supercomputing 2020

#### **Solutions (2)**

#### **Use More Memory/ Node or GPUs!**

In-efficient use of resources

-ve



#### **Out-Of-Core Methods**

- More sensitive to co-located jobs
- Tricky when no local DRAM/SSD
- Simplicity of data parallelism
- No changes to model/framework +ve
- No messing with Batch Norm (BN)

#### **Explicit Model-parallelism**

Constructing a cost model that is accurate
High implementation cost

- Analytical: deterministic

+ve

- Low overhead

### **Strategies for Distributed Training**

- ✓ Formally defines possible <u>pure parallelism strategies</u> for Distributed Deep Learning:
  - Data, Spatial, Layer, Channel, Filter, Hybrid (combine of data with others)



#### **Model Parallelism for CNN**



### **Performance and Memory Projection**

- ✓ Proposed an initial analysis/estimation model
  - Support many CNN, transformer (GEMM)
  - Based on an Ideal Parameterization

User's constraints (fixed / maximum	Parallelism	
<b>Dataset Specification</b> (Sample size, number of samples)	$D = I \times B$ samples	strategies Computation
DNN model Specification	Sumpres	Time
Network shape, Layer dimensions, e.g., $x, y, w$ . Complexity e.g. FLOP counts	$ \begin{array}{c} x[N,C, X_W \times X_H] \\ w[C,F, K \times K] \\ \cdots DV F X \times Y 1 \end{array} $	Communicati on Time
Computer System Specification Architecture breakdown	Y[N,F, T <sub>W</sub> T <sub>H</sub> ]       PLX       Switch	Max memory per PE
<ul> <li>Interconnect merarchy</li> <li>Bandwidth and latency</li> <li>Processing Element Specification</li> <li>Computation speed (FLOPS)</li> <li>Mamora consolity</li> </ul>	GPU GPU GPU GPU	Max. number of PEs

	Computation Time $T_{comp}$	Communication Time $T_{comm}$	Maximum Memory Per PE	Number of PEs $p$
Serial	$D\sum_{l=1}^{G} \left( FW_l + BW_l \right) + \frac{D}{B}\sum_{l=1}^{G} (WU_l)$	0	$2\gamma\delta\sum_{l=1}^G \left(B( x_l + y_l )+ w_l  ight)$	p = 1
Data	$\frac{D}{p}\sum_{l=1}^{G}(FW_l + BW_l) + \frac{D}{B}\sum_{l=1}^{G}(WU_l)$	$2rac{D}{B}(p-1)\Big(lpha+rac{\sum_{l=1}^G w_l }{p}\deltaeta\Big)$	$2\gamma\delta\sum_{l=1}^G \left( rac{B}{p}( x_l + y_l )+ w_l  ight)$	$p \leq B$
Spatial	$\frac{D}{p}\sum_{l=1}^{G}\left(FW_{l}+BW_{l}\right)+\frac{D}{B}\sum_{l=1}^{G}(WU_{l})$	$ \begin{array}{l} 2\frac{D}{B} \Big( (p-1)(\alpha + \frac{\sum_{l=1}^{G}  w_l }{p} \delta\beta) + \\ \sum_{l=1}^{G} \big( 2\alpha + B(halo( x_l ) + halo( \frac{dL}{dy_l} ))\delta\beta \big) \Big) \end{array} $	$2\gamma\delta\sum_{l=1}^G \left(Brac{( x_l + y_l )}{p}+ w_l  ight)$	$p = pw  imes ph \le \min_{l=1}^G (W_l  imes H_l)$
<b>Layer</b> (Pipeline)	$\frac{D(p+S-1)}{S} \left( \max_{i=1}^{p} (FW_{G_i}) + \max_{i=1}^{p} (BW_{G_i}) \right) + \max_{i=1}^{p} (WU_{G_i}) \right)$	$2\frac{D(p+S-2)}{B}\left(\max_{i=1}^{p-1}\left(\alpha+\frac{B}{S} y_{G_{i}} \delta\beta\right)\right)$	$egin{aligned} & 2\gamma\delta\max_{i=1}^p \left(\sum_{l=1}^{G_i} \left(B( x_l + y_l ) +  w_l  ight) ight) \end{aligned}$	$p \leq G$
Filter	$\frac{D}{p}\sum_{l=1}^{G} \left(FW_l + BW_l\right) + \frac{D}{Bp}\sum_{l=1}^{G} (WU_l)$	$3rac{D}{B}(p-1)\sum_{l=1}^{G-1}(lpha+rac{B y_l }{p}\deltaeta)$	$2\gamma\delta\sum_{l=1}^G \left(B( x_l + y_l )+rac{ w_l }{p} ight)$	$p \le \min_{l=1}^G (F_l)$
Channel	$\frac{D}{p}\sum_{l=1}^{G}\left(FW_{l}+BW_{l}\right)+\frac{D}{Bp}\sum_{l=1}^{G}(WU_{l})$	$3rac{D}{B}(p-1)\sum_{l=1}^{G-1}(lpha+rac{B y_l }{p}\deltaeta)$	$2\gamma\delta\sum_{l=1}^G \left(B( x_l + y_l )+rac{ w_l }{p} ight)$	$p \le \min_{l=1}^G (C_l)$
Data + Filter	$\frac{D}{p}\sum_{l=1}^{G}\left(FW_{l}+BW_{l}\right)+\frac{D}{Bp2}\sum_{l=1}^{G}(WU_{l})$	$\begin{array}{l} 3\frac{D}{B}(p2-1)\sum_{l=1}^{G-1}(\alpha+\frac{B y_l }{p}\delta\beta) + \\ 2\frac{D}{B}(p1-1)(\alpha+\frac{\sum_{l=1}^{G} w_l }{p}\delta\beta) \end{array}$	$2\gamma\delta\sum_{l=1}^{G}\left(\frac{B( x_l + y_l )}{p1}+\frac{ w_l }{p2}\right)$	$p = p1 \times p2 \le B \times \min_{l=1}^G (F_l)$

## **Improving the Estimation Accuracy?**



Estimation accuracy is not good (average ~60%), especially in communication

#### Different latency $\alpha$ and bandwidth factors $\beta \rightarrow \text{Empirical Parameterization}$

- Computation parameters  $(FW_l, BW_l, and WU_l)$
- Communication parameter ( $\alpha$  and  $\beta$ ): use NCCL-test, OSU benchmarks

#### Network contention → **Self-contention modeling**

- Introduce contention penalty coefficient  $\phi$ ,
- Using dynamic contention graphs [1] to detect contention

#### Network congestion → Detach it in the empirical result

- Report minimum communication time
- Congestion impact factor using benchmark [2]



Self-contention in hybrid parallelism with  $\phi = 2$ 

[1] Maxim Martinasso et al. 2011. A ``Contention-Aware Performance Model for HPC-Based Networks: A Case Study of the InfiniBand Network. In Euro-Par 2011 Parallel Processing."

[2] Sudheer Chunduri et al. 2019. GPCNeT: Designing a Benchmark Suite for Inducing and Measuring Contention in HPC Networks (SC '19) [HPDC] <u>Truong Thao Nguyen</u>, et al. "An Oracle for Characterizing and Guiding Large-Scale Training of Deep Neural Networks," ACM Symposium on High-Performance Parallel and Distributed Computing (HPDC2021) (to be appeared) 41

#### **Final Projection Result**



#### **Limits/Bottlenecks of Communication**



[2] N. T. Truong, M. Wahib, and R. Takano, "Efficient MPI-Allreduce for Large-Scale Deep Learning on GPU-Clusters", in Concurrency and Computation Practice and Experience, Special Issue on Parallel and Distributed Computing and Networking (CCPE)

[3]N.T. Truong, "Topology-aware Sparse Allreduce for Large-scale Deep Learning", 38th IEEE International Performance Computing and Communications Conference (IPCCC), Oct 2019

[4] N.T. Truong, et. al. "An Algorithm and System Co-design for Large-Scale Training of Distributed Deep Learning", submitted to IEEE-CCGRID 2021

### **Limits/Bottlenecks of Computation**

- Weight update
  - Weight update forms a significant portion of compute time
- Workload balancing
  - Pipeline can outperform data parallelism because of P2P rather than collective.
  - To achieve this, all stages must balance
- Computation redundancy
  - Mostly caused by frameworks
  - Convolutional layer doesn't scale as expected
  - Computation overhead in split/concat



### **Memory Capacity Limits/Bottlenecks**

- Redundancy in Memory
  - Activations are distributed but weights are not

			Memory-bound Sensitive to number of optimizer's slot variables (moving average, momentum, etc.)	
Forward pass	Backward pass	All-reduce gradients	Weight update	
			·	
SGD:				
				]
ADAM:				

- Memory Manager
  - Kernels waiting for memory leads to performance degradation (Cosmoflow – data+spatial)



# Collaboration

#### **Research Interest**





#### **Team Introduction**

We do middleware research for edge-to-cloud services and infrastructure (DigiARC calls as <u>Continuum Computing core</u>), specially we fucus on

- ① Acceleration technologies at edge, for advanced and high-performance CC-applications
- 2 Cloud technologies, which can provide huge AI computing and data processing powers with end services, through tighter integration of edge infrastructure

and contribute to the realization of various digital services considered in Society 5.0.





#### **Research Topics (add or remove?)**



#### Data management and use in CC

- Offloading and data reduction at edge for scalability and energy saving, data pipeline optimization
- DataOps automation, balanced with effective human interaction
- Data traceability in sharing and integration

#### Management of CC-application services

- Develop and build methods of high-performance and robust CC-applications for geo-distributed, unstable and heterogeneous environments
- Zero-touch service deployment, dynamic QoS control, offloading decision, mobility support
- Low-latency service-to-service interaction, autonomous service orchestration



#### High performance cloud technologies tightly connecting with edge

- Large-scale computing with accelerators
- High performance AI, AI resource hub
- Low-latency and scalable connection services with edge, efficient resource management

# TAKEAWAYS

- Large-scale training, e.g., 100-1000GPUs is challenging
  - Inference on HPC system
- New problems when training with different (Dataset, DL Model and computer System)
  - Weight update
  - I/O Stagging
  - Network congestion
- The rise of hybrid parallelism (data + model parallelism)