

#### Al For Science At Scale – Part III (Training LLMs with Hundreds of Billions of Parameters)

Sajal Dash Research Scientist, HPC and Al Analytics and Al Methods at Scale (AAIMS), NCCS dashs@ornl.gov

ORNL is managed by UT-Battelle LLC for the US Department of Energy



### A Story where Goliath Wins



2017: Attention is all you need **[Google]** 



#### 2020: Language Models are Few-Shot Learners [OpenAl]



#### Model Size Exploded with Time



**CAK RIDGE** National Laboratory

## What Do You Need to Train a Trillion Parameter Model?

Compute: 120\*P^2 ~ 120 Million ExaFlops

Memory: 20P ~ 20 Terabytes

Modern GPUs: A100 – 80GB Memory, 312Teraflops



- How to distribute the training workload across hundreds of MI250X GPUs?
- How to take advantage of the existing software developed targeting NVIDIA gpus?

## Outline

- Distributed training of LLMs
- Preparing a training software stack for Frontier
- Finding best distributed training strategies
- Training up to a Trillion parameter model on Frontier (<u>for a few</u> <u>iterations</u>)
- Training LLMs on OLCF resources





#### Forward Pass



Weights to update:  $W_K$ ,  $W_Q$ ,  $W_V$ ,  $W_1$ ,  $W_2$ 

Let, 
$$X \in \mathbb{R}^{s \times d}, W_K, W_Q, W_V \in \mathbb{R}^{d \times d}$$
.  
Then,  $K = XW_K, Q = XW_Q, V = XW_V$ ,  
 $K, Q, V \in \mathbb{R}^{s \times d}$ .  
And, Attention =  $softmax(\frac{QK^T}{\sqrt{d}})V \in \mathbb{R}^{s \times d}$   
 $\cdots$   
 $Y = func(Attention) \in \mathbb{R}^{s \times d}$   
 $W_1 \in \mathbb{R}^{d \times 4d}, W_2 \in \mathbb{R}^{4d \times d}$   
 $Z = GeLU(YW_1) \times W_2 \in \mathbb{R}^{s \times d}$ 

# Training DL Models with Large Data



Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour



**CARE RIDGE** LEADERSHIP COMPUTING FACILITY

#### Frontier Node Architecture



**CARE RIDGE** National Laboratory

#### Inside a Transformer Layer



4d^2

 $4d^{2} + 4d^{2} = 8d^{2}$ 

Number of parameters One layer: 12d^2 L layers: 12Ld^2



#### How to Grow a Transformer Model?

- Number of parameters 12Ld^2 can be think of a volume
  - With L (#Layers), the model grows linearly
  - With d (hidden dimension), model grows quadratically



-	Model	#Layers	Hidden size	#Attention heads
-	1.4B	24	2114	24
-	22B	48	6144	48
-	175B	96	12288	96
-	1T	128	25600	128



## Memory Requirement During Training an LLM

#### Model Weights:

#### Number of parameters: 12Ld^2

Memory Memory Requirement	
Optimizer States: Adam Values 22B Model 175B Model 1T Mo	del
(mea Parameters (6x) 132 GB 1050 GB 6 TH	;
• 8 bytes * number of paramet Gradients (4x) 88 GB 700 GB 4 TH	•
• 2 bytes * number of paramet Optimizer States (8x) 176 GB 1.4 TB 8 TH	•
• 4 bytes * number of paramet Total Memory (20x*) 440 GB 3.5 TB 20 T	3

• 4 bytes \* number of parameters for fp32 training

Gradients Same as number or parameters

• 4 bytes \* number of parameters for either fp32 or mixed precision training (gradients are always kept in fp32)

#### Forward Activations Batch-size x output-nodes?

• size depends on many factors, the key ones being sequence length, hidden size and batch size.



### Model Parallelism: Why and How

- Models (or Memory needed for their training) are too big to fit in a single GPU
- So, we need to break the model into pieces
- What's broken needs to be rebuilt from the pieces
- It's like Kintsugi, but with a more practical concern such as price of gold (comm. Latency)



Kintsugi



#### **Tensor Parallelism**

- Model is too large to fit in a GPU's memory
- We slice the model tensors along a suitable dimension (row or column), and the GPU memory is large enough to fit one slice.
- Unlike sharded data parallelism, this is not data parallelism, the same data gets evaluated by different part of the same layer, and the output gets combined.





#### Self-Attention





#### Feed Forward Network (FFN)





## Limitations of Tensor Parallelism

- Requires frequent AllReduce communication after every layer
- Intermediate outputs get AllReduced
- Tensor Parallel (TP) size is limited by the number of GPUs in a node (6 for Summit, 8 for Frontier)
- For TP > 6/8 the communication requires crossing node boundary through 25+25GB/s ethernet cable



#### Lessons Learnt From Tensor Parallelism





Pipeline Parallelism (PP = 2)



COAK RIDGE LEADERSHIP National Laboratory

## Pipeline Parallelism (Gpipe)







#### Pipeline Parallelism

Bubble size ~ (#Pipeline stages) / (#Microbatches)



### Pipeline Bubble vs #Microbatches

• Increasing the #Microbatches will reduce the bubble



(a) Throughput vs. global batch-size for 22B model.

(b) Throughput vs global batch-size for 1T model.

• But that will result in large global batch size, hurting the convergence



#### Bubble vs #pipeline-stages

• Reducing the #pipeline-stages reduces bubble



(a) Throughput vs. PP while keeping global batch size fixed at 128.

(b) Throughput vs. PP while scaling global batch size to keep the pipeline bubble ratio fixed.

• Then, we cannot use too many GPUs



#### Sharded Data Parallelism

- The model is too big to fit in a single GPU's memory
- One GPU has enough memory to fit a fraction of the model
- Slice the model with horizontal lines and place one horizontal slice in one GPU
- But, when a layer is being evaluated, each GPU should have the copy of that layer
- The GPU has enough memory to fit one full layer on top of its slice







## Data vs. Sharded vs. Tensor vs. Pipeline Parallelism

	Data Parallelism	Sharded Data Parallelism	Tensor Parallelism	Pipeline Parallelism
Dimension of distribution	Full Model gets replicated	Horizontal Slices + one layer replicated	Horizontal Slices	Vertical slices
Model/Output	Output	Model	Output	Output
Communicati on frequency	Least frequent	Most frequent	Most frequent	Intermediate
Communicati on volume	Large volume	Small	Small	Small
Limitation	High-volume communication	Frequent low- volume communication	Frequent low- volume communication	Bubbles from insufficient overlap
Advantage	Necessary for consuming large data	Allows large model training	Allows large model training	Can hide latency with smaller bubbles



## Outline

- Distributed training of LLMs
- Preparing a training software stack for Frontier
- Finding best distributed training strategies
- Training up to a Trillion parameter model on Frontier (<u>for a few</u> <u>iterations</u>)
- Training LLMs on OLCF resources



## Preparing AI Software Stack for Frontier

- PyTorch initially from source code and using ROCM compiler
- ROCM implementations of libraries like APEX, FlashAttention
- DeepSpeed without JIT compilation
- Utilize PyTorch's cpp\_extention and cuda\_extension for modifying cuda kernels to hip kernels



## 3D Parallelism using Megatron-DeepSpeed

- Megatron-DeepSpeed is a state-of-the-art training framework developed by Microsoft and NVIDIA
- It supports 3D+ parallelism (Tensor, Pipeline, Data, Sharded Data)
- We ported this Framework to Frontier through hipify and combining with ROCM specific packages



## Porting Megatron-DeepSpeed to Frontier

- <u>https://github.com/sajal-vt/Megatron-DeepSpeed-ORNL/tree/frontier-sd</u>
  - Install deepspeed with prebuilt ops.
     DS\_BUILD\_CPU\_ADAM=1 DS\_BUILD\_FUSED\_ADAM=1
     DS\_BUILD\_FUSED\_LAMB=1DS\_BUILD\_TRANSFORMER=1 DS\_BUILD\_AIO=1
     DS\_BUILD\_UTILS=1 pip install deepspeed
  - Following from the earlier step, commented out jit\_fusion options in megatron/training.py
     Line 34: #from megatron.initialize import set\_jit\_fusion\_options
     Line 100: #set\_jit\_fusion\_options()
  - 3. Add the following line to megatron/core/tensorparallel/layers.py line 531, config.gradient\_accumulation\_fusion = False
  - 4. In the launch script we had to set: export CUDA\_DEVICE\_MAX\_CONNECTIONS=1



## Outline

- Distributed training of LLMs
- Preparing a training software stack for Frontier
- Finding best distributed training strategies
- Training a Trillion parameter model on Frontier (<u>for a few</u> <u>iterations</u>)
- Training LLMs on OLCF resources



#### 3D Parallelism

- A Combination of Tensor, Pipeline, and Data Parallelism
- Determine how many GPUs (world-size) you need to fit the model
- Factorize world-size into TP (tensor parallel size) and PP (pipeline parallel size)

Distribution Strategy	Tunable Parameters
Tensor Parallelism	Tensor Parallel Size $(TP)$
Pipeline Parallelism	Pipeline Parallel Size (PP), #Mi-
	crobatches (m)
Sharded Data Parallelism	ZeRO-1
Common	Micro Batch Size
Mixed Precision Training	FP16, BF16

TABLE IV: Distribution Strategies and relevant tunable parameters

Hyperparameters	Range
Pipeline-parallel-size (PP)	$PP \in \{1, 2, 4, 8, 12, 16\}$
Tensor-parallel-size (TP)	$TP \in \{1, 2, 4, 8\}$
Micro-batch-size (MBS)	$MBS \in [4, 20]$
Gradient accumulation steps (GAS)	$GAS \in \{5, 10\}$
ZeRO-1 Optimizer	$ZeRO - 1 \in \{True, False\}$
Number of Nodes (NNODES)	$NNODES \in \{12, 16\}$

TABLE V: Hyperparameter Tuning for 175B Model





## Best practices with parallelism paradigms

- Tensor Parallelism
  - Keep it within the node (TP < 8)
- Pipeline Parallelism
  - Use large number of micro-batches (But that can increase the global batch-size)
- Data Parallelism
  - Can't use too much data parallelism. A large global batch size will make the model divergence.



#### Hyperparameter Search using DeepHyper



# Searching a 3D Parallelism Strategy using DeepHyper



Fig. 9: Search Trajectory of DeepHyper search



Fig. 10: Impact of various hyperparameters on training performance in terms of GPU throughput.



## Outline

- Distributed training of LLMs
- Preparing a training software stack for Frontier
- Finding best distributed training strategies
- Training a Trillion parameter model on Frontier (for a few iterations)
- Training LLMs on OLCF resources



## "Best" Strategy to Train 175B and 1T Models

Disclaimers:

- 1. We didn't train any model till completion. We only trained for 10 iterations and less than 2 hours.
- 2. We don't have any completely trained models

Hyperparameters	Value		
	175B Model	1T Model	
TP	4	8	
PP	16	64	
MBS	1	1	
GBS	640	1600	
ZeRO Stage	1	1	
Flash Attention	v2	v2	
Precision	fp16	fp16	
checkpoint-activations	True	True	

TABLE VI: Best parameters for training a 175B model and a 1T model.



Fig. 11: MI250X Throughput for various model sizes. We report the hardware FLOPS, which are in close agreement with the model FLOPS.

## Weak Scaling Performance



(a) Weak scaling of 175b model training by keeping per replica batch- (b) Weak scaling of 1T model training by keeping per replica batchsize fixed at 640.

size fixed at 1600.

Fig. 12: Weak scaling performance of 175b model and 1T model training.



## Strong Scaling Performance



(a) Strong scaling of 175b model training by keeping the total batch (b) Strong scaling of 1T model training by keeping a total batch-size size fixed at 8000. The strong scaling efficiency at 1024 GPUs is fixed at 8016. The strong scaling efficiency at 3072 GPUs is 87.05%. 89.93%.

Fig. 13: Strong scaling performance of 175b model and 1T model training.



#### Power Consumption



**COAK RIDGE** National Laboratory

43

#### Energy Budget against Model Size (~P^2)

- We used rocm-smi to measure GPU level energy consumption
- Idle GPU power is approximately 90 Watts.
- Taking iteration time, consumed tokens per iterations, average active power, and total number of MI250X GPU cards, we can estimate how much energy each model training will require.
- For 22B, 175B, and 1T parameter model, the estimated energy consumption is 284 Gigajoules, 17.65 Terajoules, and 662 Terajoules.



#### To Read About the Work...

#### **Optimizing Distributed Training on Frontier for** Large Language Models

Sajal Dash Oak Ridge National Laboratory Oak Ridge National Laboratory dashs@ornl.gov

Isaac R Lyngaas lyngaasir@ornl.gov

Junqi Yin Oak Ridge National Laboratory Oak Ridge National Laboratory yinj@ornl.gov

Xiao Wang wangx2@ornl.gov

Romain Egele Université Paris-Saclay romainegele@gmail.com

J. Austin Ellis AMD austin.ellis@amd.com

Matthias Maiterth Oak Ridge National Laboratory maiterthm@ornl.gov

**Guojing Cong** Oak Ridge National Laboratory congg@ornl.gov

Feiyi Wang Oak Ridge National Laboratory fwang2@ornl.gov

Prasanna Balaprakash Oak Ridge National Laboratory pbalapra@ornl.gov

- Arxiv Link: https://arxiv.org/abs/2312.12705
- Accepted at ISC HPC, 2024



## Outline

- Distributed training of LLMs
- Preparing a training software stack for Frontier
- Finding best distributed training strategies
- Training up to a Trillion parameter model on Frontier (<u>for a few</u> <u>iterations</u>)
- Training LLMs on OLCF resources



## Getting OLCF Allocation

<u>https://www.olcf.ornl.gov/for-users/documents-forms/olcf-directors-discretion-project-application/</u>

OLCF IN THE NEWS HPC WIRE - OLC	F'S SUMMIT SUPERCOMPUTER LIVES FOR ANOTHER YEAR A CALENDAR		
<b>CAK RIDGE</b> National Laboratory	ABOUT OLCF ▼ OLCF RESOURCES ▼ R&D ACTIVITIES ▼ SCIENCE AT OLCF ▼		
OLCF DIRECTOR'S DISCRETION PROJECT APPLICATION			
Contact Support	HOME / FOR USERS / DOCUMENTS FORMS / OLCF DIRECTORS DISCRETION PROJECT APPLICATION		
Need assistance from a trained OLCF support staff member? We're here to	The OLCF Director's Discretion Project Application Form has been moved to myOLCF.		

help.



## OLCF Tutorials on Distributed Training of Deep Learning Models

- Part I: Introduction and Data Parallel Training
  - <u>https://www.olcf.ornl.gov/calendar/ai-for-science-at-scale-intro/</u>
- Part II: Model Parallelism
  - <u>https://www.olcf.ornl.gov/calendar/ai-training-series-ai-for-science-at-</u> <u>scale-part-2/</u>



#### Hands On

- Training Large LLMs
  - Train a 22B model on 2 nodes
  - Train a 175B model on 16 nodes
  - Train a 1T model on 128 nodes
- Finding the best training strategies using DeepHyper
  - Find the best strategy for training a 22B model



# Training Large LLMs

- git clone git@github.com:sajal-vt/Megatron-DeepSpeed-ORNL.git
- git fetch
- git branch –v –a
- git switch FA2
- sbatch -reservation=ai launch\_gpt22b\_bf16.slurm
- sbatch --reservation=ai launch\_gpt175b\_bf16.slurm
- sbatch --reservation=ai launch\_gpt1T\_bf16.slurm



# Finding the best training strategies using DeepHyper

- git clone git@github.com:sajal-vt/Megatron-DeepSpeed-ORNL.git
- git fetch
- git branch –v –a
- git switch frontier-sd
- sbatch --reservation=ai launch\_dh.frontier

