### Introduction to Parallel Computing (CMSC416 / CMSC616)





### Parallel Deep Learning

### Abhinav Bhatele, Alan Sussman

### Annoucements

- Assignment 7 (extra credit) has been posted
  - Due on Dec 9 11:59 pm (no extensions for any reason)
- IRB-approved research study to analyze the generated logs
  - Please opt-in to help us with research



Abhinav Bhatele (CMSC416 / CMSC616)

2

### The evolution of HPC systems and rise of a new revolution in Al

- In the last two decades, an enormous amount of compute power has become available
- Large datasets and open source software such as PyTorch have also emerged
- Led to a frenzy in the world of AI and the effects are being felt in almost every other domain









### The evolution of HPC systems and rise of a new revolution in Al

- In the last two decades, an enormous amount of compute power has become available
- Large datasets and open source software such as PyTorch have also emerged
- Led to a frenzy in the world of AI and the effects are being felt in almost every other domain









# The evolution of HPC systems and rise of a new revolution in Al

- In the last two decades, an enormous amount of compute power has become available
- Large datasets and open source software such as PyTorch have also emerged
- Led to a frenzy in the world of AI and the effects are being felt in almost every other domain





### Use of AI/ML in HPC

- approaches to model performance and power
- 2018-present: Deep learning and transfer learning starting to become popular



• Some of the earliest works were around 2007 on using traditional machine learning

• 2007-2017: Supervised and unsupervised algorithms for creating prediction models



### Use of AI/ML in HPC

- approaches to model performance and power
- 2018-present: Deep learning and transfer learning starting to become popular



Abhinav Bhatele, et al. Identifying the culprits behind network congestion. In Proceedings of the IEEE International Parallel & Distributed Processing Symposium, IPDPS '15. May 2015.



• Some of the earliest works were around 2007 on using traditional machine learning

2007-2017: Supervised and unsupervised algorithms for creating prediction models



### Use of AI/ML in HPC

- approaches to model performance and power
- 2018-present: Deep learning and transfer learning starting to become popular



Abhinav Bhatele, et al. Identifying the culprits behind network congestion. In Proceedings of the IEEE International Parallel & Distributed Processing Symposium, IPDPS '15. May 2015.



• Some of the earliest works were around 2007 on using traditional machine learning

2007-2017: Supervised and unsupervised algorithms for creating prediction models



Daniel Nichols, et al. Predicting Cross-Architecture Performance of Parallel Programs. In Proceedings of the IEEE International Parallel & Distributed Processing Symposium. May 2024.





## **AI/ML for Computational Science**

- Using of AI models for approximating computation in scientific codes
- Black box models
- Physics-informed machine learning mod



d	e	S
Ч		J



## **AI/ML for Computational Science**

- Using of AI models for approximating computation in scientific codes
- Black box models
- Physics-informed machine learning mod





d	e	S





## **AI/ML for Computational Science**

- Using of AI models for approximating computation in scientific codes
- Black box models
- Physics-informed machine learning models









### Deep neural networks (DNNs)

- An area of machine learning that uses artificial neural networks to learn complex functions
  - Often from high-dimensional data: text, images, audio, ...
- Widespread use in computer vision, natural language processing, etc.
- Neural networks can be used to model complex functions
- Several layers that process "batches" of the input data







# Deep neural networks (DNNs)

- An area of machine learning that uses artificial neural networks to learn complex functions
  - Often from high-dimensional data: text, images, audio, ...
- Widespread use in computer vision, natural language processing, etc.
- Neural networks can be used to model complex functions
- Several layers that process "batches" of the input data









# Deep neural networks (DNNs)

- An area of machine learning that uses artificial neural networks to learn complex functions
  - Often from high-dimensional data: text, images, audio, ...
- Widespread use in computer vision, natural language processing, etc.
- Neural networks can be used to model complex functions
- Several layers that process "batches" of the input adata

X<sub>m</sub>



Summation Activation Inputs Weights and bias function

Outputs





Abhinav Bhatele @ PEARC<sup>24</sup>





https://medium.com/data-science-at-microsoft/how-large-language-models-work-91c362f5b78f





Abhinav Bhatele @ HPDC '24













Word		Probability
speak		0.065
gener	ate	0.072
politic	s	0.001
walk		0.003

Abhinav Bhatele @ HPDC '24









Word	Probability		Word	Probability
speak	0.065		ability	0.002
generate	0.072	7	text	0.084
politics	0.001		coherent	0.085
walk	0.003		ideas	0.041



### **Other definitions**

- Learning/training: task of selecting weights that lead to an accurate function
- Loss: a scalar proxy that when minimized leads to higher accuracy
- Gradient descent: process of updating the weights using gradients (derivates) of the loss weighted by a learning rate
- Mini-batch: Small subsets of the dataset processed iteratively
- Epoch: One pass over all the mini-batches





- The largest model you can run on an H100 96 GB GPU is around 3.5-4 billion parameters
- On a single node (with four H100 GPUs): around ~16 billion parameters model
- Training a 16B parameter would take 33 years!
- OpenAl's GPT 4.0 is estimated to have 1.8 trillion parameters
- Meta's Llama-3.1-405B has more than 400 billion parameters



### Do we really need parallel resources?









### Why is LLM training well-suited for HPC?











### **Parallel/distributed training**

- Many opportunities for exploiting paral
- Iterative process of training (epochs)
- Many iterations per epoch (mini-batche
- Many layers in DNNs



Abhinav Bhatele (CMSC416 / CMSC616)

llelism	
es)	

### **Parallel/distributed training**

- Many opportunities for exploiting parallelism
- Iterative process of training (epochs)
- Many iterations per epoch (mini-batches)
- Many layers in DNNs



Abhinav Bhatele (CMSC416 / CMSC616)



### Parallel/distributed training

- Many opportunities for exploiting parallelism
- Iterative process of training (epochs)
- Many iterations per epoch (mini-batches)
- Many layers in DNNs





2048 GPUs

3072 CPUs

400 GPUs

KARMA

LBANN

ZeRO

Data

Data

Data

17B

78.6B

100B

### Sequential LLM training

while (remaining batches) { Read a single batch

PSSG

Forward pass: perform matrix multiplies to compute output activations, and a loss on the batch

Backward pass: matrix multiplies to compute gradients o the loss w.r.t. parameters via backpropagation

Optimizer step: use gradients to update the weights or parameters such that loss is gradually reduced





f		
4		

- Divide training data (input batch) among workers (GPUs)
- Each worker has a full copy of the entire NN and processes different mini-batches
- All reduce operation to synchronize gradients
- Example: PyTorch's DDP, ZeRO







- Divide training data (input batch) amor workers (GPUs)
- Each worker has a full copy of the entitient NN and processes different mini-batch
- All reduce operation to synchronize gradients
- Example: PyTorch's DDP, ZeRO



g				
e				
es	Batch			

Abhinav Bhatele @ PEARC<sup>24</sup>



- Divide training data (input batch) amor workers (GPUs)
- Each worker has a full copy of the entitient NN and processes different mini-batch
- All reduce operation to synchronize gradients
- Example: PyTorch's DDP, ZeRO



	Shard 2	
s Batch	Shard I	
	Shard 0	



- Divide training data (input batch) among workers (GPUs)
- Each worker has a full copy of the entire NN and processes different mini-batches
- All reduce operation to synchronize gradients
- Example: PyTorch's DDP, ZeRO





Abhinav Bhatele @ PEARC<sup>24</sup>



- Divide training data (input batch) among workers (GPUs)
- Each worker has a full copy of the entire NN and processes different mini-batches
- All reduce operation to synchronize gradients
- Example: PyTorch's DDP, ZeRO





Abhinav Bhatele @ PEARC<sup>24</sup>



### Inter-layer parallelism

- Assign entire layers to different processes/GPUs
  - Ideally map contiguous subsets of layers
- managing different layers
- Use a pipeline of mini-batches to enable concurrent execution











### Point-to-point communication (activations and gradients) between processes/GPUs



### Inter-layer parallelism

- Assign entire layers to different processes/GPUs
  - Ideally map contiguous subsets of layers
- managing different layers
- Use a pipeline of mini-batches to enable concurrent execution







Point-to-point communication (activations and gradients) between processes/GPUs



### Intra-layer parallelism

- Enables training neural networks that would not fit on a single GPU
- Distribute the work within each layer to multiple processes/GPUs
  - Essentially parallelize matrix operations such as matmuls across multiple GPUs
- Example: Megatron-LM







### Intra-layer parallelism

- Enables training neural networks that would not fit on a single GPU
- Distribute the work within each layer to multiple processes/GPUs
  - Essentially parallelize matrix operations such as matmuls across multiple GPUs
- Example: Megatron-LM





Tensor parallelism



# Hybrid parallelism

- Using two or more approaches together in the same parallel framework
- 3D parallelism: use all three
- Popular serial frameworks: pytorch, tensorflow
- Popular parallel frameworks: DDP, MeshTensorFlow, Megatron-LM, ZeRO



Abhinav Bhatele (CMSC416 / CMSC616)



• A hybrid parallelism approach







• A hybrid parallelism approach

Batch







• A hybrid parallelism approach

Batch







• A hybrid parallelism approach

Batch







• A hybrid parallelism approach

Batch

 Combines data parallelism with **3-dimensional parallel matrix** multiplication (PMM)





### Data Parallelism





### Enabling 3D parallel matrix multiplication in AxoNN

- Distribute I and W across a 3D grid of GPUs
- Compute partial output activations, O on each GPU











# AxoNN

- Distribute I and W across a 3D grid of GPUs
- Compute partial output activations, O on each GPU



















# AxoNN

- Distribute I and W across a 3D grid of GPUs
- Compute partial output activations, O on each GPU















# AxoNN

- Distribute I and W across a 3D grid of GPUs
- Compute partial output activations, O on each GPU















# **5** Easy parallelization using AxoNN

- Requires minimal code changes to model architecture (code):

  - with auto parallelize():
- Our ML collaborators used this mode for the memorization experiments
- We also have backends for lightning and accelerate

### PSSG

from axonn.intra layer import auto parallelize

net = # declare your sequential model here

AxoNN intercepts all declarations of torch.nn.Linear, and parallelizes them













Singh & Bhatele @









Singh & Bhatele @





### PSSG



Singh & Bhatele @































# UNIVERSITY OF MARYLAND

### Questions?

