

CS 598 Al Efficiency: Systems and Algorithms Overview & Key Challenges in LLMs Training Systems

Minjia Zhang

Computer Science Department

The Large Language Model Revolution

What is Language Model?



Evolution of DNN Models

15,000x increase in 5 years



Model size

Dialogue/New Search



ChatGPT: Optimizing Language Models for Dialogue

Code Continuation and Generation

Python 3 Get the current value of a Bitcoin in US dollars using the bitcoincharts api """
<pre>import requests import json</pre>
<pre>def get_bitcoin_price(): url = 'http://api.bitcoincharts.com/v1/weighted_prices.json' response = requests.get(url) data = json.loads(response.text) return data['USD']['7d']</pre>
<pre>ifname == 'main': print(get_bitcoin_price())</pre>

<u>Suggest code and entire function in your editor – Github/OpenAl Codex</u>

Image Generation from Text

TEXT DESCRIPTION

An astronaut Teddy bears A bowl of soup

that is a portal to another dimension that looks like a monster as a planet in the universe

as digital art in the style of Basquiat drawn on a cave wall

 \rightarrow

DALL-E 2





DALL·E: Creating Images from Text - OpenAl

Multi-Agent LLM Applications



AutoGen: Enabling Next-Gen LLM Applications via Multi-Agent Conversation

Autonomous Driving



A Language Agent for Autonomous Driving

Transformers for Language Modeling



Attention Is All You Need, NeurIPS 2017



BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, ACL 2019

Language Models are Few-Shot Learners, NeurIPS 2020

LLMs are Impressively Scaling!



Figure 2 We show a series of language model training runs, with models ranging in size from 10^3 to 10^9 parameters (excluding embeddings).

Scaling Laws for Neural Language Models, OpenAl, 2020

AI Efficiency Challenges

- Too slow to train high-quality models on massive data
 - More hardware ≠ higher throughput, bigger model
 - Higher throughput ≠ better accuracy, faster convergence, lower cost
 - Better techniques ≠ handy to use
- Slow and expensive to deploy the models

DL System Desired Capabilities (3E)

Efficiency: Efficient use of hardware for high scalability and throughput

Effectiveness: High accuracy and fast convergence, lowering cost

Easy to use: Improve development productivity of model scientists

ML/DL Training Problem Definition Recap

- Given model f, data set $\{xi, y_i\}_{i=1}^N$
- Minimize the loss between predicted labels and true labels: $Min \frac{1}{N} \sum_{i=1}^{N} loss(f(x_i, y_i))$
- Common loss function
 - Cross-entropy, MSE (mean squared error)
- Common way to solve the minimization problem
 - Stochastic gradient descent (SGD)
 - Adaptive learning rates optimizers (e.g., Adam)

Gradient Descent

- Model f_w is parameterized by weight w
- $\eta > 0$ is the learning rate



Adaptive Learning Rates (Adam)

- Model f_w is parameterized by weight w
- $\eta > 0$ is the learning rate

For t = 1 to T

$$\Delta w = \eta x \frac{1}{N} \sum_{i=1}^{N} \nabla \left(loss(f_w(x_i, y_i)) \right)$$

$$w = \Delta w // \text{ apply update}$$
End

$$u_t = \beta_1 * \nu_{t-1} - (1 - \beta_1) * g_t$$

 $s_t = \beta_2 * s_{t-1} - (1 - \beta_2) * g_t^2$

 $\Delta \omega_t = -\eta \frac{\nu_t}{\sqrt{s_t + \epsilon}} * g_t$

 $g_t : Gradient at time t along \omega^j$

 $\nu_t : Exponential Average of gradients along \omega_j$

 $s_t : Exponential Average of squares of gradients along \omega_j$

 $g_t : Gradients along \omega_j$

Adam: A Method for Stochastic Optimization, 2014

Accelerating Gradient Descent

- Model f_w is parameterized by weight w
- $\eta > 0$ is the learning rate

For t = 1 to T $\Delta w = \eta \times \frac{1}{N} \sum_{i=1}^{N} \nabla \left(loss(f_w(x_i, y_i)) \right) // \text{ compute derivative and update}$ $w \to \Delta w // \text{ apply update}$ End

Data Parallelism (DP)



Implemented as standard component in DL training frameworks, such as PyTorch DDP

1. Partition the training data

2. Parallel training on different machines

3. Synchronize the local updates

4. Refresh local model with new parameters, then go to 2

Scaling Distributed Machine Learning with the Parameter Server, 2014

Distributed Data Parallel Training in GPU Clusters



Training Efficiency: Breaking the Memory Wall

Large Model Training Challenges

	Bert-		Turing	
	Large	GPT-2	17.2 NLG	GPT-3
Parameters	0.32B	1.5B	17.2B	175B
Layers	24	48	78	96
Hidden Dimension	1024	1600	4256	12288
Relative Computation	1x	4.7x	54x	547x
Memory Footprint	5.12GB	24GB	275GB	2800GB



Large Model Training Challenges

	Bert-		Turing	
	Large	GPT-2	17.2 NLG	GPT-3
Parameters	0.32B	1.5B	17.2B	175B
Layers	24	48	78	96
Hidden Dimension	1024	1600	4256	12288
Relative Computation	1x	4.7x	54x	547x
Memory Footprint	5.12GB	24GB	275GB	2800GB

Out of Memory

NVIDIA V100 GPU memory capacity: 16G/32G NVIDIA A100 GPU memory capacity: 40G/80G



DNN Training Hits the Memory Wall



*AI and Memory Wall. (This blogpost has been written in... | by Amir Gholami | riselab | Medium

Distributed Training Techniques

- Tensor Parallelism
- Pipeline Parallelism
- 3D Parallelism
- ZeRO-Style Data Parallelism



Tensor Parallelism

Splice tensors across GPUs

+

synchronization primitives
(e.g., all-reduce)





Supported in:

- DeepSpeed
- Megatron-LM

Megatron-LM: Training Multi-Billion Parameter Language Models Using Model Parallelism, 2019

Pipeline Parallelism



- **PyTorch**
- **DeepSpeed**
- **Megatron-LM**

- Naïve model parallelism leads to severe underutilization
- Gpipe divides batch into micro-batches, enabling different device to work on different micro-batches, reducing pipeline bubbles and improving utilization

GPipe: Efficient Training of Giant Neural Networks using Pipeline Parallelism, 2019

3D Parallelism



DeepSpeed Extreme Scale Model Training For Everyone, 2020

ZeRO-Style Data Parallelism

- ZeRO removes the redundancy across data parallel process
- Partitioning optimizer states, gradients and parameters (3 stages)



Supported in: •<u>DeepSpeed</u> •PyTorch

Large Models Need Parallelism

	Max Parameter (in billions)	Max Parallelism	Compute Efficiency	Usability (Model Rewrite)
Data Parallel (DP)	Approx. 1.2	>1000	Very Good	Great
Tensor Parallel (TP)	Approx. 20	Approx. 16	Good	Needs Model Rewrite
TP + DP	Approx. 20	> 1000	Good	Needs Model Rewrite
Pipeline Parallel (PP)	Approx. 100	Approx. 128	Very Good	Needs Model Rewrite
PP + DP	Approx. 100	> 1000	Very Good	Needs Model Rewrite
TP + PP + DP	> 1000	> 1000	Very Good	Needs Significant Model Rewrite
ZeRO	> 1000	> 1000	Very Good	Great

Long Context LLMs

Foundation Model Context Length



Variable Sequence Length Training for Long-Context Large Language Model, 2023

Long-Context Training Systems



Reducing Activation Recomputation in Large Transformer Models [Arxiv 2022] DeepSpeed Ulysses: System Optimizations for Enabling Training of Extreme Long Sequence Transformer Models [PODC 2024]



Gradient Checkpointing/Rematerialization



Step 2) Divide the network into segments before backpropagation



A Possible Allocation Plan



Training Deep Nets with Sublinear Memory Cost, 2016



Coop: Memory is not a Commodity, 2023 ³²

New Hardware Support



Introducing the NVIDIA H100 Tensor Core GPU [2022]

Mixed Precision Training



Mixed Precision Training [ICLR, 2018]

FP8-LM: Training FP8 Large Language Models [Arxiv, 2023]



Figure 8: Training losses of GPT-125M models with the settings presented in Tab. 6. The loss curve for FP8 #4 has diverged.

Memory/Storage Hierarchy

Heterogeneous Memory





ZeRO-Offload: Democratizing Billion-Scale Model Training [ATC, 2021] ZeRO-Infinity: Breaking the GPU Memory Wall for Extreme Scale Deep Learning [SC, 2021]



Auto-Parallelism for Easy-to-Use



(a) Computational graph



(c) The space of intra-operator parallelism (e.g., Tofu [55])



(d) The space of inter-operator parallelism (e.g., DAPPLE [17])



(b) Manual plan (e.g., Megatron-LM [40])



(e) Our hierarchical space

Alpa: Automating Inter- and Intra-Operator Parallelism for Distributed Deep Learning [OSDI, 2022]

Proteus: Simulating the Performance of Distributed DNN Training [Arxiv, 2023]



QA