Memory Footprint Reduction Techniques for DNN Training: An Overview

Episode III. Weight Placement for EX-Large Models

Bojian Zheng 2022/2/18

Outline

- EX-Large Models: Bring us better accuracy but also challenges.
- Existing Solutions: Data and Model Parallelism
- ZeRO-Infinity: Data-Parallelism-based Solution for EX-Large Training

Recall: A History of Prior Works

...

Weight Pruning for Efficient **Inference**

2015-2016

Feature Maps Reduction for Efficient **Training**

2016-

Weight Placement for **EX-Large Models**

Topic for Today!

...

2019-

- [1] T. Chen et al. *DianNao: A Small-Footprint High-Throughput Accelerator for Ubiquitous Machine-Learning*. ASPLOS 2014
- [2] S. Han et al. *EIE: Efficient Inference Engine on Compressed Deep Neural Network*. ISCA 2015
- [3] S. Han et al. *Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding*. ICLR 2015
- [4] Y. Chen et al. Eyeriss: A Spatial Architecture for Energy-Efficient Dataflow for Convolutional Neural Networks. ISCA 2016

- [5] T. Chen et al. *Training Deep Nets with Sublinear Memory Cost.* arXiv 2016
- [6] M. Rhu et al. vDNN: Virtualized Deep Neural Networks for Scalable, Memory-Efficient Neural Network Design. MICRO 2016
- [7] A. Jain et al. *Gist: Efficient Data Encoding for Deep* Neural Network Training. ISCA 2018
- [8] B. Zheng et al. Echo: Compiler-based GPU Memory Footprint Reduction for LSTM RNN Training. ISCA 2020

[9] S. Rajbhandari et al. *ZeRO-Infinity: Breaking the GPU Memory Wall for Extreme Scale Deep Learning.* SC 2021

Recall: GPU Memory Allocations in DNN Training



- Major GPU memory consumers: Model States & Feature Maps
 - Model States: weights (1), gradients (2), optimizer states (3)

Why EX-Large Models?

• Rule of Thumb in Deep Learning:

Larger & Deeper Models \Rightarrow **Better Model Accuracy**^[1, 2]

• The amount of computation resources consumed by a job can reject the **problem scale** and may also indicate the **commercial value** of the workload^[3].

[1] K. He et al. Deep Residual Learning for Image Recognition. CVPR 2016

[2] T. Brown et al. Language Models are Few-Shot Learners. NeurIPS 2020

[3] M. Wang et al. Characterizing Deep Learning Training Workloads on Alibaba-PAI. IISWC 2019

Challenges with Training EX-Large Models

• BERT^[1]

- State-of-the-Art Natural Language Processing Model (2018)
- 340M parameters
- Cannot fit into the GPU memory even with a batch size of 1^[2].
- Addressable with offloading^[3], checkpointing^[4] etc.



- [1] J. Devlin et al. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. NAACL 2019
- [2] https://github.com/google-research/bert
- [3] M. Rhu et al. vDNN: Virtualized Deep Neural Networks for Scalable, Memory-Efficient Neural Network Design. MICRO 2016
- [4] T. Chen et al. Training Deep Nets with Sublinear Memory Cost. arXiv 2016

Challenges with Training EX-Large Models

• GPT-3^[1]

- State-of-the-Art Natural Language Processing Model (**2020**)
- 175B parameters (470× LARGER)
- 700 GB storage ≫ modern GPU memory capacity
- 3×700 GB = 2.1 TB more for gradients & optimizer states



Challenges with Training EX-Large Models

• Multi-Interests^[1, 2, 3]

- Industry-scale recommender
- Encodes every object of interest into a vector, which assembles into an **embedding table**.
- Up to 240 GB storage ≫ modern GPU memory capacity
- Cannot fit a single layer into the GPU memory!
- [1] P. Covington et al. *Deep Neural Networks for YouTube Recommendations*. RecSys 2016
- [2] J. Weston et al. *Nonlinear Latent Factorization by Embedding Multiple User Interests.* RecSys 2013
- [3] M. Wang et al. Characterizing Deep Learning Training Workloads on Alibaba-PAI. IISWC 2019



Today's DNN Training In A Nutshell



Data and Model Parallelism

• Partition into multiple devices? E.g., Data and/or Model Parallelism?



 $\bullet \bullet \bullet$

Data Parallelism

• Partitions the **data** into multiple devices.



Model Parallelism

• Partitions the **model** into multiple devices.



 $\bullet \bullet \bullet$

Data vs. Model Parallelism

Data Parallelism

+

╋

Model Parallelism

- Need to significantly refactor the code for load-balancing.
- Hard to support models with complex dependencies.

 Need to replicate the entire model. Challenges of the EXlarge models persist!

00

Could we have the merits from both? i.e., generic support for EX-Large models with minor code changes?

+



- A solution that is based on data-parallelism.
- Key Idea Partitions model states among devices and use efficient communication primitives to gather them.
- However, hungry DNN models need much more memory!



ZeRO-Infinity^[1]

- ZeRO^[2] + Infinity Offload Engine
- Observation GPUs are also accompanied with a CPU and NVMe SSDs

	GPUs	CPU	NVMe SSDs
Memory (TB)	0.5	1.5	28

 Key Idea Partitions model states among (GPUs + CPU + SSD) and Keeps most of the states to the latter two.



• Challenge #1 How to efficiently get the weights from CPU/SSDs?

• Challenge #2 How to handle cases when the weight of a single layer cannot fit into the GPU memory?

- Challenge #1 How to efficiently get the weights from CPU/SSDs?
- Optimization #1 ∀weight, each GPU worker holds a portion of it and gathers the rest from others.
 - Better utilization than having a single worker broadcast the weight to others.

		Aggregate	Aggregate	Aggregate			CPU-GPU	NVME-GPU
		GPU	CPU	NVME	GPU Memory	GPU-GPU	Memory	Memory
		Memory	Memory	Memory	Bandwidth/GPU	Bandwidth	Bandwidth/GPU	Bandwidth/GPU
Nodes	GPUs	(TBs)	(TBs)	(TBs)	(GB/s)	(GB/s)	(GB/s)	(GB/s)
1	1	0.032	1.5	28.0	600-900	N/A	12.0	12.0
1	16	0.5	1.5	28.0	600-900	150-200	4.1	1.6
4	64	2.0	6.0	112.0	600-900	60-100	4.1	1.6
16	256	8.0	24.0	448.0	600-900	60-100	4.1	1.6
64	1024	32.0	96.0	1792.0	600-900	60-100	4.1	1.6

- Challenge #1 How to efficiently get the weights from CPU/SSDs?
- Optimization #1 ∀weight, each GPU worker holds a portion of it and gathers the rest from others.
 - Better utilization than having a single worker broadcast the weight to others.

		Aggregate	Aggregate	Aggregate			CPU-GPU	NVME-GPU
		GPU	CPU	NVME	GPU Memory	GPU-GPU	Memory	Memory
		Memory	Memory	Memory	Bandwidth/GPU	Bandwidth	Bandwidth/GPU	Bandwidth/GPU
Nodes	GPUs	(TBs)	(TBs)	(TBs)	(GB/s)	(GB/s)	(GB/s)	(GB/s)
1	1	0.032	1.5	28.0	600-900	N/A	12.0	12.0
1	16	0.5	1.5	28.0	600-900	150-200	4.1	1.6
4	64	2.0	6.0	112.0	600-900	60-100	4.1	1.6
16	256	8.0	24.0	448.0	600-900	60-100	4.1	1.6
64	1024	32.0	96.0	1792.0	600-900	60-100	4.1	1.6
1 4 16 64	16 64 256 1024	0.5 2.0 8.0 32.0	1.5 6.0 24.0 96.0	28.0 112.0 448.0 1792.0	600-900 600-900 600-900 600-900	150-200 60-100 60-100 60-100	4.1 4.1 4.1 4.1	1.6 1.6 1.6 1.6

- Challenge #1 How to efficiently get the weights from CPU/SSDs?
- Optimization #1 ∀weight, each GPU worker holds a portion of it and gathers the rest from others.
 - Better utilization than having a single worker broadcast the weight to others.

		Aggregate	Aggregate	Aggregate			CPU-GPU	NVME-GPU
		GPU	CPU	NVME	GPU Memory	GPU-GPU	Memory	Memory
		Memory	Memory	Memory	Bandwidth/GPU	Bandwidth	Bandwidth/GPU	Bandwidth/GPU
Nodes	GPUs	(TBs)	(TBs)	(TBs)	(GB/s)	(GB/s)	(GB/s)	(GB/s)
1	1	0.032	1.5	28.0	600-900	N/A	12.0	12.0
1	16	0.5	1.5	28.0	600-900	150-200	4.1	1.6
4	64	2.0	6.0	112.0	600-900	60-100	4.1	1.6
16	256	8.0	24.0	448.0	600-900	60-100	4.1	1.6
64	1024	32.0	96.0	1792.0	600-900	60-100	4.1	1.6

- Challenge #1 How to efficiently get the weights from CPU/SSDs?
- Optimization #1 ∀weight, each GPU worker holds a portion of it and gathers the rest from others.
 - Better utilization than having a single worker broadcast the weight to others.

		Aggregate	Aggregate	Aggregate			CPU-GPU	NVME-GPU
		GPU	CPU	NVIVIE	GPU Memory	GPU-GPU Bondwidth	Memory Rendwidth (CDU	Nemory Rendwidth (CDU
Nodes	GPUs	(TBs)	(TBs)	(TBs)	(GB/s)	(GB/s)	(GB/s)	(GB/s)
1	1	0.032	1.5	28.0	600-900	N/A	12.0	12.0
1	16	0.5	1.5	28.0	600-900	150-200	4.1	1.6
4	64	2.0	6.0	112.0	600-900	60-100	4.1	1.6
16	256	8.0	24.0	448.0	600-900	60-100	4.1	1.6
64	1024	32.0	96.0	1792.0	600-900	60-100	4.1	1.6

Same effective bandwidth even if 1024 GPUs are reading in parallel

- Challenge #1 How to efficiently get the weights from CPU/SSDs?
- Optimization #1 ∀weight, each GPU worker holds a portion of it and gathers the rest from others.
 - Better utilization than having a single worker broadcast the weight to others.

		Aggregate	Aggregate	Aggregate			CPU-GPU	NVME-GPU
		GPU	CPU	NVME	GPU Memor	GPU-GPU	Memory	Memory
		Memory	Memory	Memory	Bandwidth/GP	U Bandwidth	Bandwidth/GPU	Bandwidth/GPU
Nodes	GPUs	(TBs)	(TBs)	(TBs)	(GB/s)	(GB/s)	(GB/s)	(GB/s)
1	1	0.032	1.5	28.0	600-900	N/A	12.0	12.0
1	16	0.5	1.5	28.0	600-900	150-200	4.1	1.6
4	64	2.0	6.0	112.0	600-900	60-100	4.1	1.6
16	256	8.0	24.0	448.0	600-900	60-100	4.1	1.6
64	1024	32.0	96.0	1792.0	600-900	60-100	4.1	1.6

4/1.6 TB/s (Gather) vs. 12/12 GB/s (Broadcast)

Key Idea #1. Pipelined Prefetch

- Challenge #1 How to efficiently get the weights from CPU/SSDs?
- Optimization #1 ∀weight, each GPU worker holds a portion of it and gathers the rest from others.

• Optimization #2 Pipelined prefetch from SSD \rightarrow CPU \rightarrow GPU.

i	SSD	CPU	GPU	Compute
<i>i</i> + 1		SSD	CPU	GPU
<i>i</i> + 2			SSD	CPU
<i>i</i> + 3				SSD

Read the weight slice from SSD for layer i + 3 while layer i is computing its results.

Key Idea #3. Tiled Computation

- Challenge #2 How to handle cases when the weight of a single layer cannot fit into the GPU memory?
- Key Idea Some layers can evaluate their results in tiles and some do not need the whole weight for the full results^[1].
 - E.g., Dense Layer $Y = XW^T$





Baseline with Data + Model Parallelism



Offloading to CPU/SSDs cause 2%

performance degradation





27

Conclusions

- EX-Large Models: Bring us better accuracy but also challenges.
- Existing Solutions: Data and Model Parallelism
 - Can't make EX-large training practical while being easily programmable.
- ZeRO-Infinity: Data-Parallelism-based Solution for EX-Large Training
 - Key Ideas: Weight Gathering, Pipelined Prefetch, Tiled Computation
 - Makes it possible to train models with trillion parameters with small performance overhead (2%).
- <u>https://github.com/microsoft/DeepSpeed</u>

Memory Footprint Reduction Techniques for DNN Training: An Overview

Episode III. Weight Placement for EX-Large Models

Bojian Zheng 2022/2/18

Why Not?

- Reduce the precision of the weight parameters^[1]? E.g., FP32→FP16?
- Challenge: Weights need to be in high-precision format to accommodate for incremental gradient updates.

• E.g., Type equation here.





ZeRO-Infinity^[1]

- ZeRO^[2] + Infinity Offload Engine
- Observation GPUs are accompanied with large CPU memory and
- Key Idea Partitions model states among workers (GPUs + CPU + NVMe SSD). Offload the states in the latter two



[1] S. Rajbhandari et al. ZeRO: Memory Optimizations toward Training Trillion Parameter Models. SC 2020



