# Scaling Deep Learning

CS246: Mining Massive Datasets
Jure Leskovec, Stanford University
Charilaos Kanatsoulis, Stanford University
http://cs246.stanford.edu

# Announcements

- **Colab 9 is due today (3/13 at 11:59 pm)**
- **Final Review Session**
  - **Friday at 8 pm via zoom (see ed for the link)**
- **Final on 3/20 at 12:15 pm**
  - **Hewlett Packard 200**
  - **2 cheat sheets allowed (front and back)**
  - **Exam score: All content found in Lecs/Slides/ Homework/ Colabs**
    - **Focus of the exam: Lectures 2 – 19**
    - **No coding**
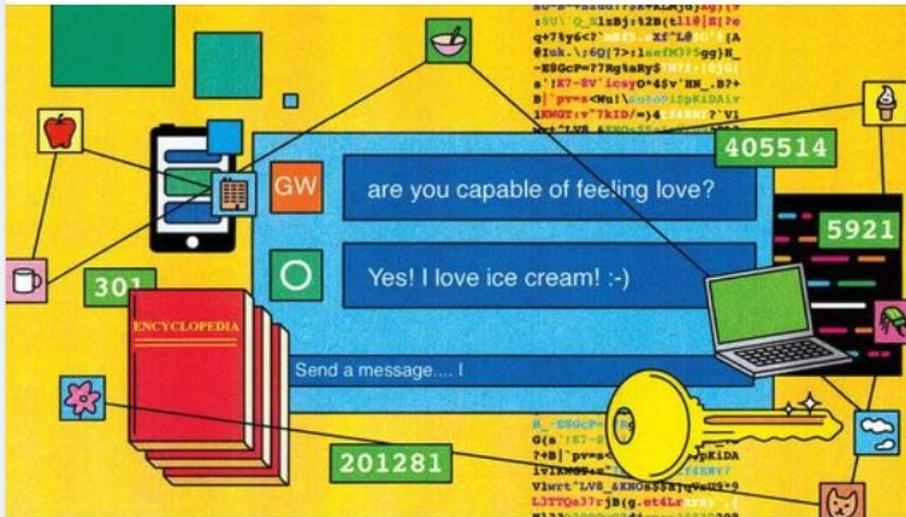  - **See ed for more details**

# Large Language Models

are you capable of feeling love?

Yes! I love ice cream! :-)

Send a message.... I

**Generative AI**

## How does ChatGPT actually work?

*Despite the feeling of magic, large language models (LLMs) are, in reality, a giant exercise in statistics*

## Language Models are Few-Shot Learners

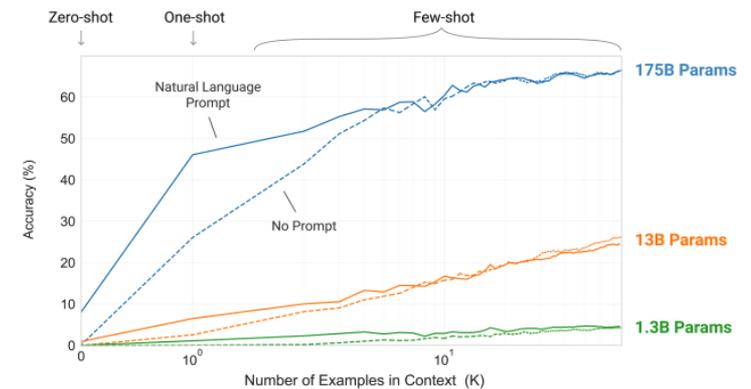| Tom B. Brown* | Benjamin Mann* | Nick Ryder* | Melanie Subbiah* |
| Jared Kaplan† | Prafulla Dhariwal | Arvind Neelakantan | Pranav Shyam | Girish Sastry |
| Amanda Askell | Sandhini Agarwal | Ariel Herbert-Voss | Gretchen Krueger | Tom Henighan |
| Rewon Child | Aditya Ramesh | Daniel M. Ziegler | Jeffrey Wu | Clemens Winter |
| Christopher Hesse | Mark Chen | Eric Sigler | Mateusz Litwin | Scott Gray |
| | Benjamin Chess | Jack Clark | Christopher Berner | |
| Sam McCandlish | Alec Radford | Ilya Sutskever | Dario Amodei |

OpenAI

**Figure 1.2: Larger models make increasingly efficient use of in-context information.** We show in-context learning performance on a simple task requiring the model to remove random symbols from a word, both with and without a natural language task description (see Sec. 3.9.2). The steeper "in-context learning curves" for large models demonstrate improved ability to learn a task from contextual information. We see qualitatively similar behavior across a wide range of tasks.

# Foundation Models for Cell Biology

# Foundation Models for Cell Biology

- **UCE creates universal representations of cells**

**Input:**
RNA expression of a single cell/nucleus

**Output:**
Cell Embedding

# Foundation Models for Cell Biology

## Visualize and transfer annotations



## Infer Hierarchies

# Training Large Models

- **Modern AI trains very large models with a huge amount of data**

- **GPT-3 is trained with 500B tokens, more than 600 GB of Data**

- **The model has 175B parameters, requiring 350 GB of storage space (<span style="color:red">GPT-4 at 1.8T parameters</span>)**

- **The memory capacity of modern GPUs is 24-192 GB**

- **There is need for developing large-scale methods that can train such models**

# Today's Lecture

- **How to train large deep learning models?**

- **Memory Optimization Methods**

- **Parallel and distributed training with multiple GPUs**
  - **Model Parallel Training**
  - **Data Parallel Training**

# Memory Optimization

# GPU Architecture hierarchy

| SM | SM | SM |
|----|----|----|
| Shared Memory | L1 | Shared Memory | L1 | Shared Memory | L1 |

L2 cache memory

Global memory

- **RTX 3090: 7GPU clusters, 84 SMs per cluster, 24 GB memory**

- **H100: 8 GPU clusters, 144 SMs per cluster, 80 GB memory**

# Sources of Memory Consumption

- **Input Data:** sequencies, images, graphs, etc.

- **Trainable parameters**

- **Auxiliary optimization variables**

- **Intermediate activation values and gradients**

- **Training Deep Nets with Sublinear Memory Cost [Chen et al., 2016]**
  - https://arxiv.org/pdf/1604.06174.pdf

# Computation graph

Input    Linear    ReLU    Linear    Softmax    Output

# Computing Gradients

| Input | Linear | ReLU | Linear | Softmax | Loss |
|-------|--------|------|--------|---------|------|
| □ → | □ → | □ → | □ → | □ → | □ |
| ↓ | ↓ | ↓ | ↓ | ↓ | ↑ |
| □ ← | □ ← | □ ← | □ ← | □ ← | □ |
| Linear-grad | ReLU-grad | Linear-grad | Softmax-grad | Loss-grad | Label |

# Checkpointing

Input      Linear      ReLU      Linear      Softmax      Output

# Checkpointing with backprop

Input　　　Linear　　　ReLU　　　Linear　　　Softmax　　　Loss

Linear-grad　　ReLU-grad　　Linear-grad　　Softmax-grad　　Loss-grad　　Label

| Input | Linear | ReLU | Linear | Softmax | Loss |
|-------|--------|------|--------|---------|------|

| Linear-grad | ReLU-grad | Linear-grad | Softmax-grad | Loss-grad | Label |
|-------------|-----------|-------------|--------------|-----------|-------|

# Checkpointing with backprop

| Input | Linear | ReLU | Linear | Softmax | Loss |
|-------|--------|------|--------|---------|------|

| | | | | | |
|--|--|--|--|--|--|

| Linear-grad | ReLU-grad | Linear-grad | Softmax-grad | Loss-grad | Label |

# Checkpointing with backprop

| Input | Linear | ReLU | Linear | Softmax | Loss |
|-------|--------|------|--------|---------|------|

| Linear-grad | ReLU-grad | Linear-grad | Softmax-grad | Loss-grad | Label |
|-------------|-----------|-------------|--------------|-----------|-------|

# Checkpointing with backprop

Input  Linear  ReLU  Linear  Softmax  Loss

Linear-grad  ReLU-grad  Linear-grad  Softmax-grad  Loss-grad  Label

# Sublinear Memory Cost

- The memory cost of the previous approach is:
  $O\left(\frac{N}{2}\right)$ for N Neural Network Layers and
  $O(1)$ for additional computations.
- If we checkpoint every K layers, the total memory cost is:

$$O\left(\frac{N}{K}\right) + O(K)$$

- **For $K = \sqrt{N}$ we reach sublinear memory cost!**

# Sublinear Memory Cost



(a) Feature map memory cost estimation

(b) Runtime total memory cost

Source: Training Deep Nets with Sublinear Memory Cost [Chen et al., 2016]

(a) ResNet

Source: Training Deep Nets with Sublinear Memory Cost [Chen et al., 2016]

# Model Parallel Training

# Model Parallelism

- **What if the model does not fit into GPU memory?**

- Idea: Split the model into *submodels* and fit each *submodel* into a different GPU

# Model Parallel Training



Layer 1    Layer 2    Layer 3

Loss

GPU1    GPU2    GPU3

# Parallel Training



- Move input data GPU1
  - Run forward pass
- Move activations from GPU1 to GPU2
  - Run forward pass
- Move activations from GPU2 to GPU3
  - Run forward pass
  - Compute loss
  - Run backward pass
- Move gradients from GPU3 to GPU2
  - Run backward pass
- Move gradients form GPU2 to GPU1
  - Run backward pass
  - Apply gradient descent step

# Challenges

- **How to move activations and gradients?**
  - **Via CPU: bad idea**
  - **Move things between GPUs**

- **How to reconcile for dependencies?**
  - **Pipelining**

# Model Parallel Training

**Figure 2:** An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input is 150,528-dimensional, and the number of neurons in the network's remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–4096–1000.

Source: Imagenet classification with deep convolutional neural networks
[Krizhevsky et al., 2012]

# Model parallel training recap

- **Works for models that do not fit GPU memory**

- **Requires pipeline design to reconcile for dependencies**

- Pipelining requires synchronization
  - Not always that easy

# Data Parallel Training

# Data Parallelism

- **What if the model fits into GPU memory, but minibatches do not fit?**

- **Use ideas from HPC**

- Idea: Split minibatches into smaller batches and feed each of them into a different GPU

# Data Parallelism

- **For each minibatch we need to compute**

$$\theta = \theta - \frac{\alpha}{n} \sum_{i=1}^{n} \nabla \mathcal{L}(f_\theta(x_i), y_i)$$

- Split each minibatch to K smaller batches

$$\theta_k = \theta - \frac{\alpha}{n} \sum_{i=1}^{n/K} \nabla \mathcal{L}(f_\theta(x_i), y_i)$$

$$\theta = \frac{1}{K} \sum_{k=1}^{K} \theta_k$$

- First idea: Use the AllReduce framework

# AllReduce



- Each GPU communicates d derivatives to K-1 GPUs

# AllReduce

- **Total communication cost per GPU is**

  **d (K-1)**

- Linear scaling with the number of GPUs

- Communication becomes a bottleneck

# Parameter Server

GPU1

GPU2

GPU4 (server)

GPU3

- Each GPU communicates d derivates to 1 server
- Server broadcasts d derivatives to K workers:
- **Communication cost:** d K for the server

# Parallel Parameter Server



GPU1

GPU2

GPU3

GPU4 (server)

GPU4 (server)

- **Each GPU** communicates d derivatives to 1 server
- **Server broadcasts** d derivatives to K / M workers
- **Communication cost:** d K / M for server

# Breaking the Limits of Param. Server

- **Can we do better?**

- Idea: Decompose the **all-reduce** operations into separate **reduce-scatter** and **all-gather** operations

# Ring-AllReduce



Source: https://engineering.fb.com/2021/07/15/open-source/fsdp/

# Ring-AllReduce

## Phase 1: Reduce-scatter
- We divide the array in each GPU into chunks
- The gradients corresponding to the same chunk index are sequentially summed across all GPUs
- Each GPU has a fully aggregated gradient for one chunk

## Phase 2: All-gather
- The fully aggregated gradients on each GPU are made available to all GPUs.

# Ring-AllReduce

Arrays Being Summed

| GPU 0 | $a_0$ | $b_0$ | $c_0$ | $d_0$ | $e_0$ |
|-------|-------|-------|-------|-------|-------|

| GPU 1 | $a_1$ | $b_1$ | $c_1$ | $d_1$ | $e_1$ |
|-------|-------|-------|-------|-------|-------|

| GPU 2 | $a_2$ | $b_2$ | $c_2$ | $d_2$ | $e_2$ |
|-------|-------|-------|-------|-------|-------|

| GPU 3 | $a_3$ | $b_3$ | $c_3$ | $d_3$ | $e_3$ |
|-------|-------|-------|-------|-------|-------|

| GPU 4 | $a_4$ | $b_4$ | $c_4$ | $d_4$ | $e_4$ |
|-------|-------|-------|-------|-------|-------|

Source: https://engineering.fb.com/2021/07/15/open-source/fsdp/

# Ring-AllReduce



Source: https://engineering.fb.com/2021/07/15/open-source/fsdp/

# Reduce-Scatter

Arrays Being Summed

| GPU 0 | $a_0$ | $b_0$ | $c_0$ | $d_0$ | $e_0$ |
|-------|-------|-------|-------|-------|-------|
| GPU 1 | $a_1$ | $b_1$ | $c_1$ | $d_1$ | $e_1$ |
| GPU 2 | $a_2$ | $b_2$ | $c_2$ | $d_2$ | $e_2$ |
| GPU 3 | $a_3$ | $b_3$ | $c_3$ | $d_3$ | $e_3$ |
| GPU 4 | $a_4$ | $b_4$ | $c_4$ | $d_4$ | $e_4$ |

Source: https://engineering.fb.com/2021/07/15/open-source/fsdp/

# Reduce-Scatter

Arrays Being Summed

| GPU 0 | $a_0$ | $b_0$ | $c_0$ | $d_0$ | $e_0+e_4$ |
|---|---|---|---|---|---|
| GPU 1 | $a_1+a_0$ | $b_1$ | $c_1$ | $d_1$ | $e_1$ |
| GPU 2 | $a_2$ | $b_2+b_1$ | $c_2$ | $d_2$ | $e_2$ |
| GPU 3 | $a_3$ | $b_3$ | $c_3+c_2$ | $d_3$ | $e_3$ |
| GPU 4 | $a_4$ | $b_4$ | $c_4$ | $d_4+d_3$ | $e_4$ |

Source: https://engineering.fb.com/2021/07/15/open-source/fsdp/

# Reduce-Scatter

| | | | | |
|---|---|---|---|---|
| **GPU 0** | $a_0$ | $b_0$ | $c_0$ | $d_4+d_3+d_0$ | $e_0+e_4$ |
| **GPU 1** | $a_1+a_0$ | $b_1$ | $c_1$ | $d_1$ | $e_0+e_4+e_1$ |
| **GPU 2** | $a_1+a_0+a_2$ | $b_2+b_1$ | $c_2$ | $d_2$ | $e_2$ |
| **GPU 3** | $a_3$ | $b_2+b_1+b_3$ | $c_3+c_2$ | $d_3$ | $e_3$ |
| **GPU 4** | $a_4$ | $b_4$ | $c_3+c_2+c_4$ | $d_4+d_3$ | $e_4$ |

Source: https://engineering.fb.com/2021/07/15/open-source/fsdp/

# Reduce-Scatter

| GPU 0 | $a_0$ | $b_0$ | $c_3+c_2+c_4+c_0$ | $d_4+d_3+d_0$ | $e_0+e_4$ |
|---|---|---|---|---|---|
| GPU 1 | $a_1+a_0$ | $b_1$ | $c_1$ | $d_4+d_3+d_0+d_1$ | $e_0+e_4+e_1$ |
| GPU 2 | $a_1+a_0+a_2$ | $b_2+b_1$ | $c_2$ | $d_2$ | $e_0+e_4+e_1+e_2$ |
| GPU 3 | $a_1+a_0+a_2+a_3$ | $b_2+b_1+b_3$ | $c_3+c_2$ | $d_3$ | $e_3$ |
| GPU 4 | $a_4$ | $b_2+b_1+b_3+b_4$ | $c_3+c_2+c_4$ | $d_4+d_3$ | $e_4$ |

Source: https://engineering.fb.com/2021/07/15/open-source/fsdp/

# Reduce-Scatter

| GPU 0 | $a_0$ | $b_2+b_1+b_3+b_4+b_0$ | $c_3+c_2+c_4+c_0$ | $d_4+d_3+d_0$ | $e_0+e_4$ |
|---|---|---|---|---|---|

| GPU 1 | $a_1+a_0$ | $b_1$ | $c_3+c_2+c_4+c_0+c_1$ | $d_4+d_3+d_0+d_1$ | $e_0+e_4+e_1$ |
|---|---|---|---|---|---|

| GPU 2 | $a_1+a_0+a_2$ | $b_2+b_1$ | $c_2$ | $d_4+d_3+d_0+d_1+d_2$ | $e_0+e_4+e_1+e_2$ |
|---|---|---|---|---|---|

| GPU 3 | $a_1+a_0+a_2+a_3$ | $b_2+b_1+b_3$ | $c_3+c_2$ | $d_3$ | $e_0+e_4+e_1+e_2+e_3$ |
|---|---|---|---|---|---|

| GPU 4 | $a_1+a_0+a_2+a_3+a_4$ | $b_2+b_1+b_3+b_4$ | $c_3+c_2+c_4$ | $d_4+d_3$ | $e_4$ |
|---|---|---|---|---|---|

Source: https://engineering.fb.com/2021/07/15/open-source/fsdp/

# All-gather

| | | | | | |
|---|---|---|---|---|---|
| **GPU 0** | $a_0$ | $b_2+b_1+b_3+b_4+b_0$ | $c_3+c_2+c_4+c_0$ | $d_4+d_3+d_0$ | $e_0+e_4$ |
| **GPU 1** | $a_1+a_0$ | $b_1$ | $c_3+c_2+c_4+c_0+c_1$ | $d_4+d_3+d_0+d_1$ | $e_0+e_4+e_1$ |
| **GPU 2** | $a_1+a_0+a_2$ | $b_2+b_1$ | $c_2$ | $d_4+d_3+d_0+d_1+d_2$ | $e_0+e_4+e_1+e_2$ |
| **GPU 3** | $a_1+a_0+a_2+a_3$ | $b_2+b_1+b_3$ | $c_3+c_2$ | $d_3$ | $e_0+e_4+e_1+e_2+e_3$ |
| **GPU 4** | $a_1+a_0+a_2+a_3+a_4$ | $b_2+b_1+b_3+b_4$ | $c_3+c_2+c_4$ | $d_4+d_3$ | $e_4$ |

# All-gather



| GPU 0 | $a_1+a_0+a_2+a_3+a_4$ | $b_2+b_1+b_3+b_4+b_0$ | $c_3+c_2+c_4+c_0$ | $d_4+d_3+d_0$ | $e_0+e_4$ |
| GPU 1 | $a_1+a_0$ | $b_2+b_1+b_3+b_4+b_0$ | $c_3+c_2+c_4+c_0+c_1$ | $d_4+d_3+d_0+d_1$ | $e_0+e_4+e_1$ |
| GPU 2 | $a_1+a_0+a_2$ | $b_2+b_1$ | $c_3+c_2+c_4+c_0+c_1$ | $d_4+d_3+d_0+d_1+d_2$ | $e_0+e_4+e_1+e_2$ |
| GPU 3 | $a_1+a_0+a_2+a_3$ | $b_2+b_1+b_3$ | $c_3+c_2$ | $d_4+d_3+d_0+d_1+d_2$ | $e_0+e_4+e_1+e_2+e_3$ |
| GPU 4 | $a_1+a_0+a_2+a_3+a_4$ | $b_2+b_1+b_3+b_4$ | $c_3+c_2+c_4$ | $d_4+d_3$ | $e_0+e_4+e_1+e_2+e_3$ |

Source: https://engineering.fb.com/2021/07/15/open-source/fsdp/

# All-gather

| | | | | |
|---|---|---|---|---|
| **GPU 0** | $a_1+a_0+a_2+a_3+a_4$ | $b_2+b_1+b_3+b_4+b_0$ | $c_3+c_2+c_4+c_0+c_1$ | $d_4+d_3+d_0+d_1+d_2$ | $e_0+e_4+e_1+e_2+e_3$ |
| **GPU 1** | $a_1+a_0+a_2+a_3+a_4$ | $b_2+b_1+b_3+b_4+b_0$ | $c_3+c_2+c_4+c_0+c_1$ | $d_4+d_3+d_0+d_1+d_2$ | $e_0+e_4+e_1+e_2+e_3$ |
| **GPU 2** | $a_1+a_0+a_2+a_3+a_4$ | $b_2+b_1+b_3+b_4+b_0$ | $c_3+c_2+c_4+c_0+c_1$ | $d_4+d_3+d_0+d_1+d_2$ | $e_0+e_4+e_1+e_2+e_3$ |
| **GPU 3** | $a_1+a_0+a_2+a_3+a_4$ | $b_2+b_1+b_3+b_4+b_0$ | $c_3+c_2+c_4+c_0+c_1$ | $d_4+d_3+d_0+d_1+d_2$ | $e_0+e_4+e_1+e_2+e_3$ |
| **GPU 4** | $a_1+a_0+a_2+a_3+a_4$ | $b_2+b_1+b_3+b_4+b_0$ | $c_3+c_2+c_4+c_0+c_1$ | $d_4+d_3+d_0+d_1+d_2$ | $e_0+e_4+e_1+e_2+e_3$ |

Source: https://engineering.fb.com/2021/07/15/open-source/fsdp/
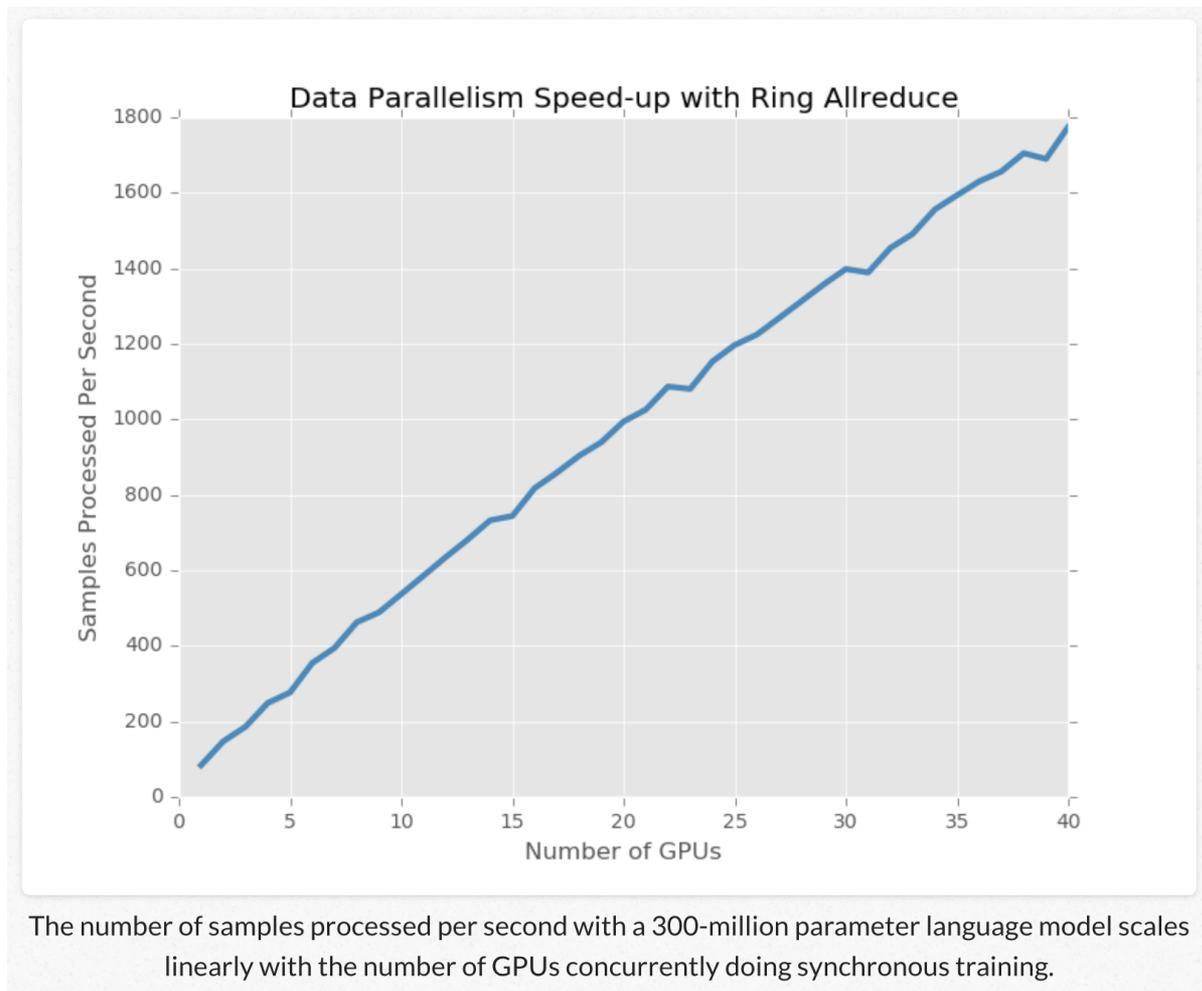
# Cost of Ring-AllReduce

- Each GPU sends and receives values K-1 times for reduce-scatter, and K-1 times for the all-gather.

- Each time, the GPUs will send d / K values

- The total cost for every GPU is $2d\,\dfrac{K-1}{K}$

# Ring-AllReduce in practice



Data Parallelism Speed-up with Ring Allreduce

The number of samples processed per second with a 300-million parameter language model scales linearly with the number of GPUs concurrently doing synchronous training.

Source: https://engineering.fb.com/2021/07/15/open-source/fsdp/

# Data parallelism recap

- **AllReduce**
  - Cost per GPU: $d\,(K-1)$

- **Parameter Server**
  - Cost per GPU: $d$
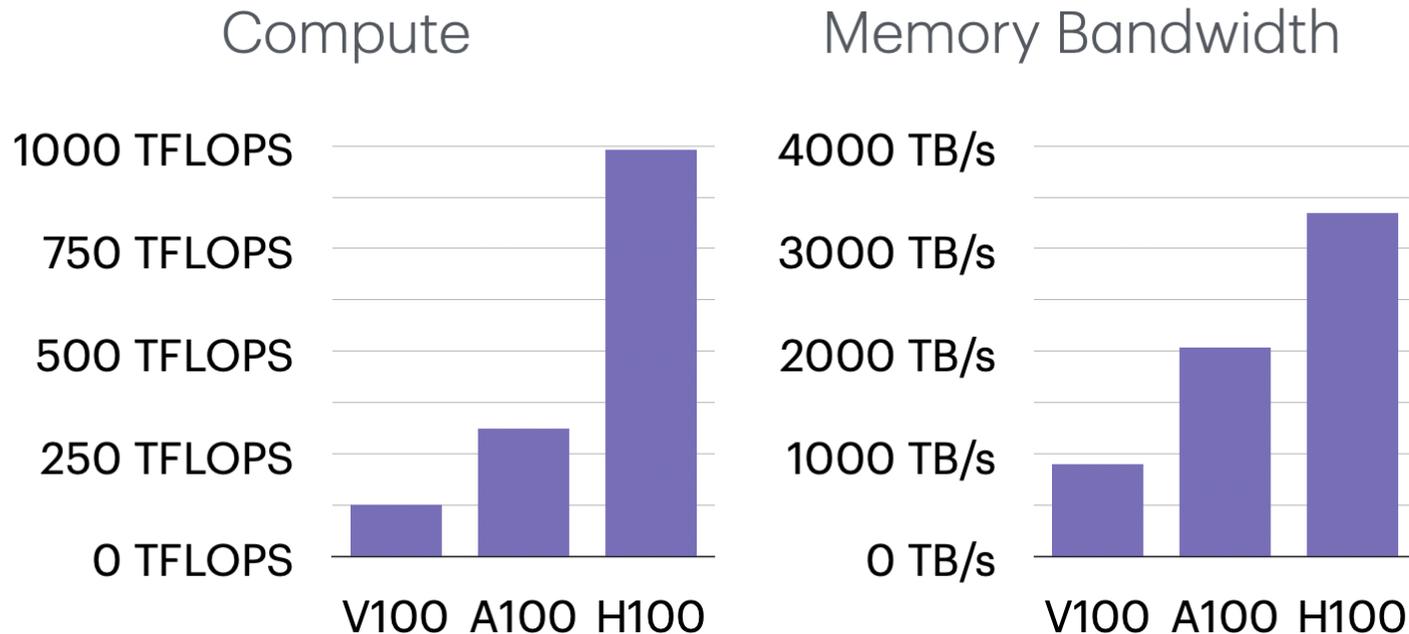  - Cost per Server: $d\,(K-1)/M$

- **Ring-AllReduce**
  - Cost per GPU: $2d\,\dfrac{K-1}{K}$

# Memory Bandwidth scaling

# Memory access patterns

- Compute scaling has outpaced memory bandwidth

Compute                    Memory Bandwidth

| Compute | |
|---|---|
| 1000 TFLOPS | |
| 750 TFLOPS | |
| 500 TFLOPS | |
| 250 TFLOPS | |
| 0 TFLOPS | |
| V100 A100 H100 | |

| Memory Bandwidth | |
|---|---|
| 4000 TB/s | |
| 3000 TB/s | |
| 2000 TB/s | |
| 1000 TB/s | |
| 0 TB/s | |
| V100 A100 H100 | |

- Memory access patterns and memory bandwidth optimization matter a lot