# Reducing Deep Learning training times with parallelism strategies and TensorFlow Distribute.

Juan Manuel Muñoz Betancur
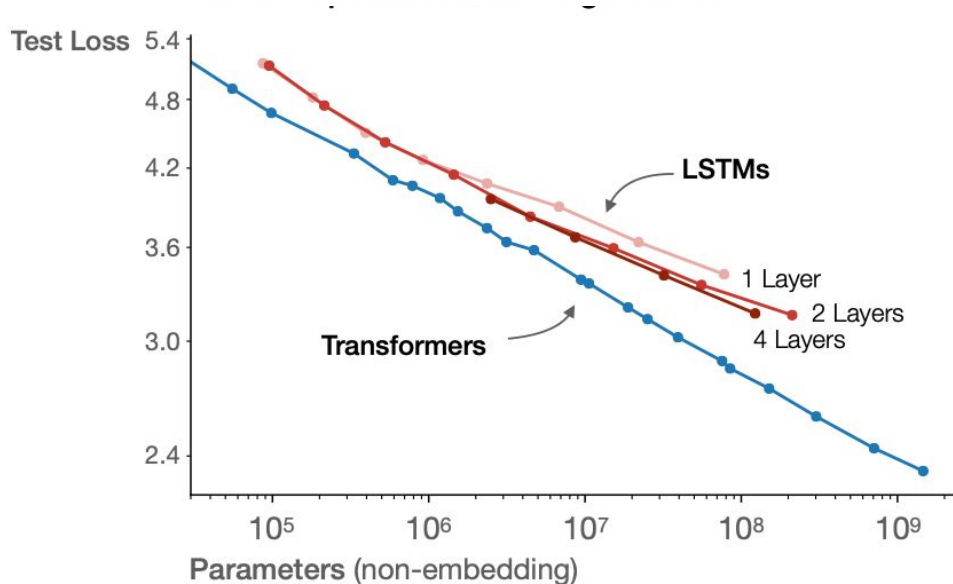
# PRESENTATION OUTLINE

- Introduction/Motivation.
- Common types of parallelism.
  - Model Parallelism
  - Data Parallelism
- Tensorflow Demo.

# Motivation

Get results faster and iterate quickly in your small/medium/large projects without additional costs.

State of the art models: GPT-3 has 175 Billion parameters, to train it even in the best SINGLE GPU available in the market would take about 355 years.

# Motivation

To train a model you need to take into account the size of the parameters, the forward pass activations (saving them for the backprop) and the gradients of the weights.

**Params size**
InceptionV3 24M params
@ fp64 = 24 x 10^6 x 8 bytes
**= 192 MB**

**Activations size**
InceptionV3 @ fp64 batch_size=256
~= 256 (4 x 10^6 x 8 bytes )
**= 32MB * 256 = 8GB**

**Params size**
GPT-3 175B parameters @fp16
= 175 x 10^9 x 2 bytes
**= 316.2 GB**

**Activations size**
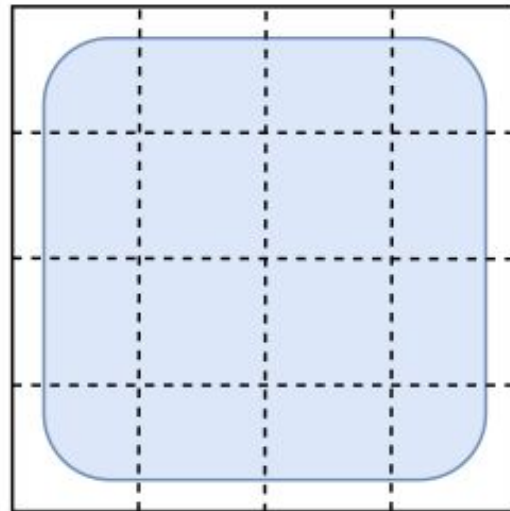GPT-3 **~TBs**

# COMMON TYPES OF PARALLELISM

- Model Parallelism.
  - Pipeline Parallelism
- Data Parallelism.
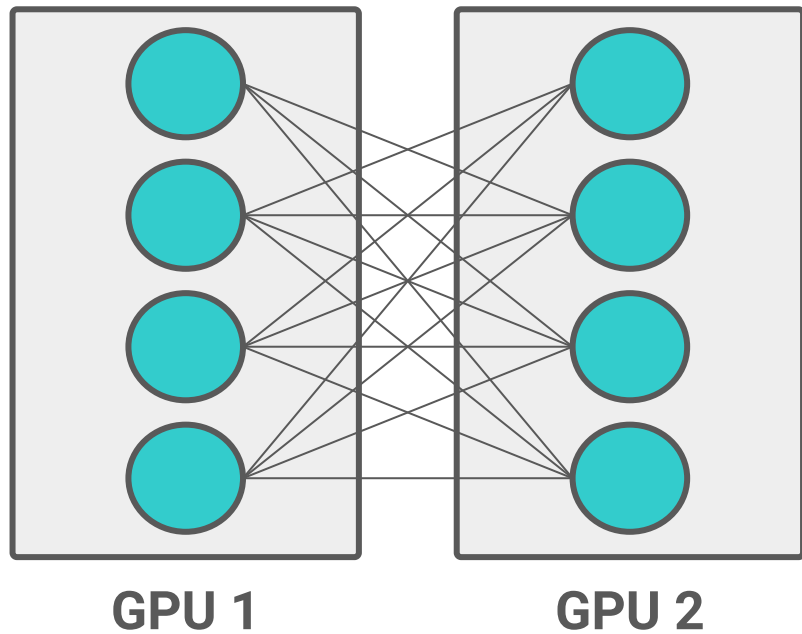
# Model Parallelism (Large model training)

**split the model across devices**

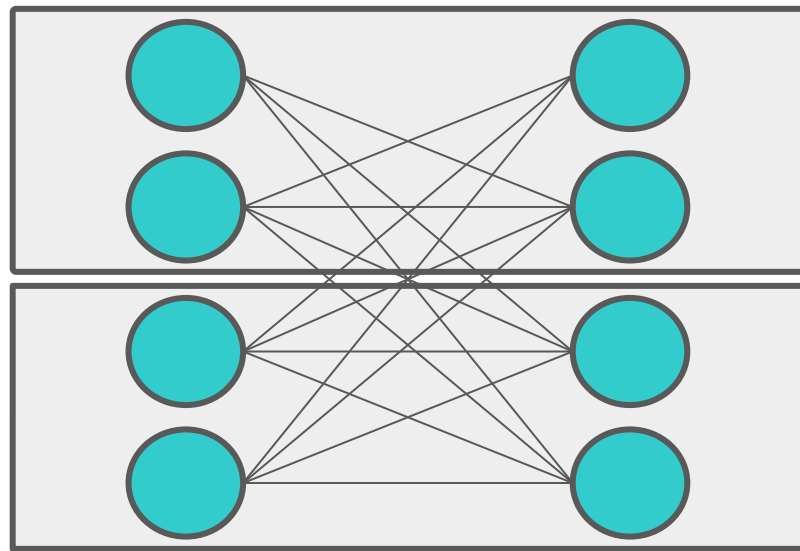each device runs a fragment of the model
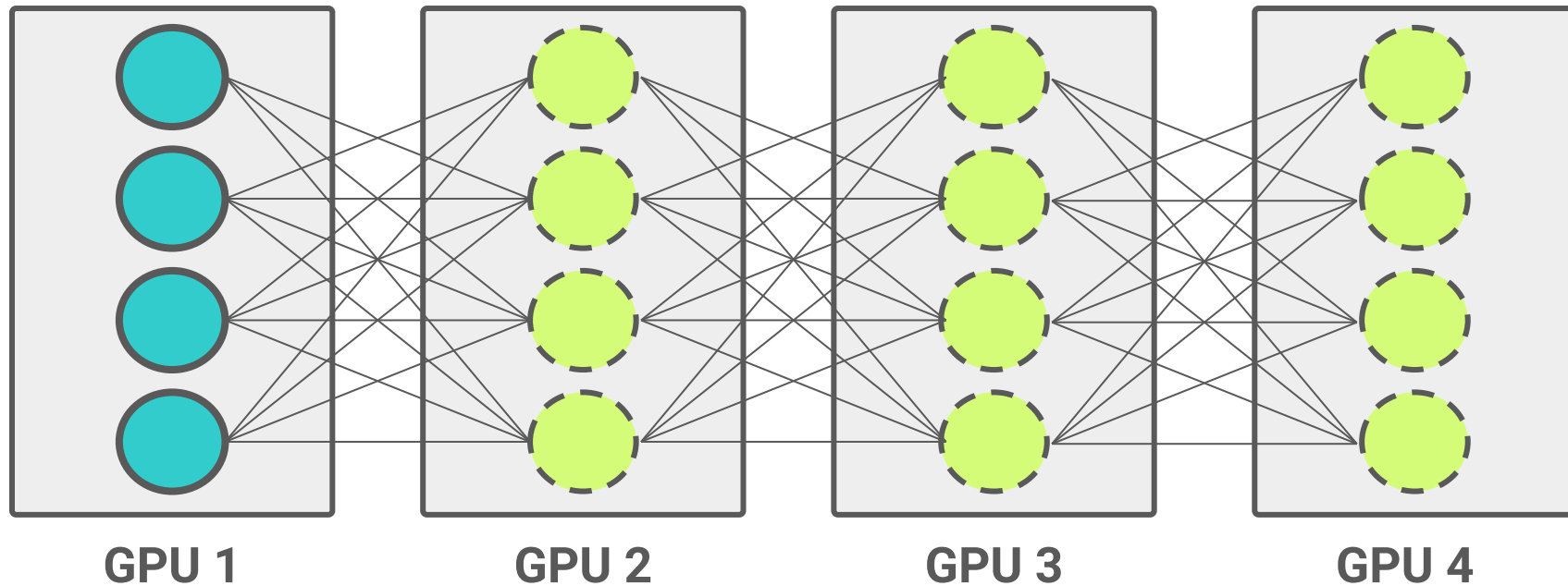
# Pipeline Parallelism

# Distributed Tensor Computation



**GPU 1**
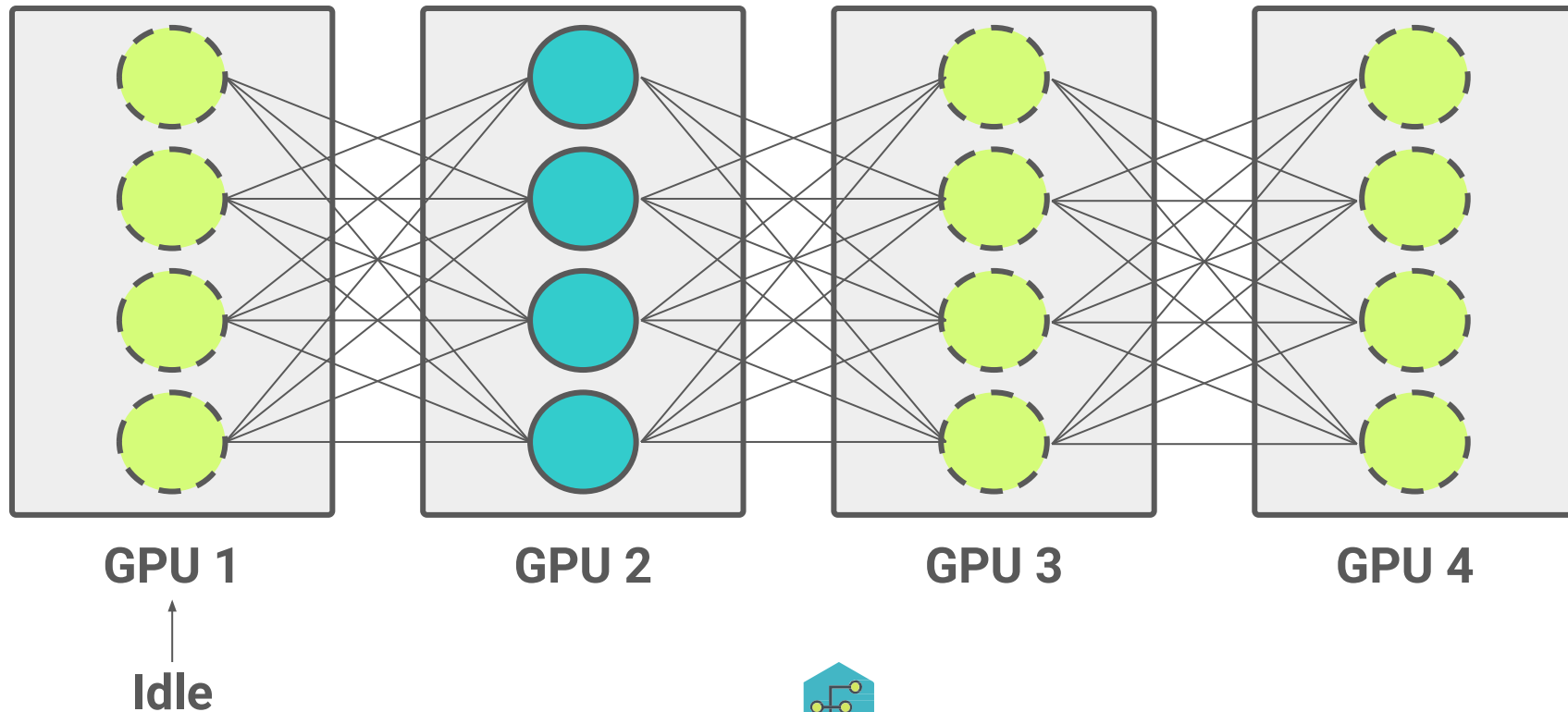
**GPU 2**

**GPU 1**

**GPU 2**

Credit: Scaling Laws for Neural Language Models, stanford cs329 slides

7

# Model Parallelism: Naive



**GPU 1**  **GPU 2**  **GPU 3**  **GPU 4**

Credit: Scaling Laws for Neural Language Models, stanford cs329 slides

# Model Parallelism: Naive



GPU 1          GPU 2          GPU 3          GPU 4

Idle

# Model Parallelism: Naive



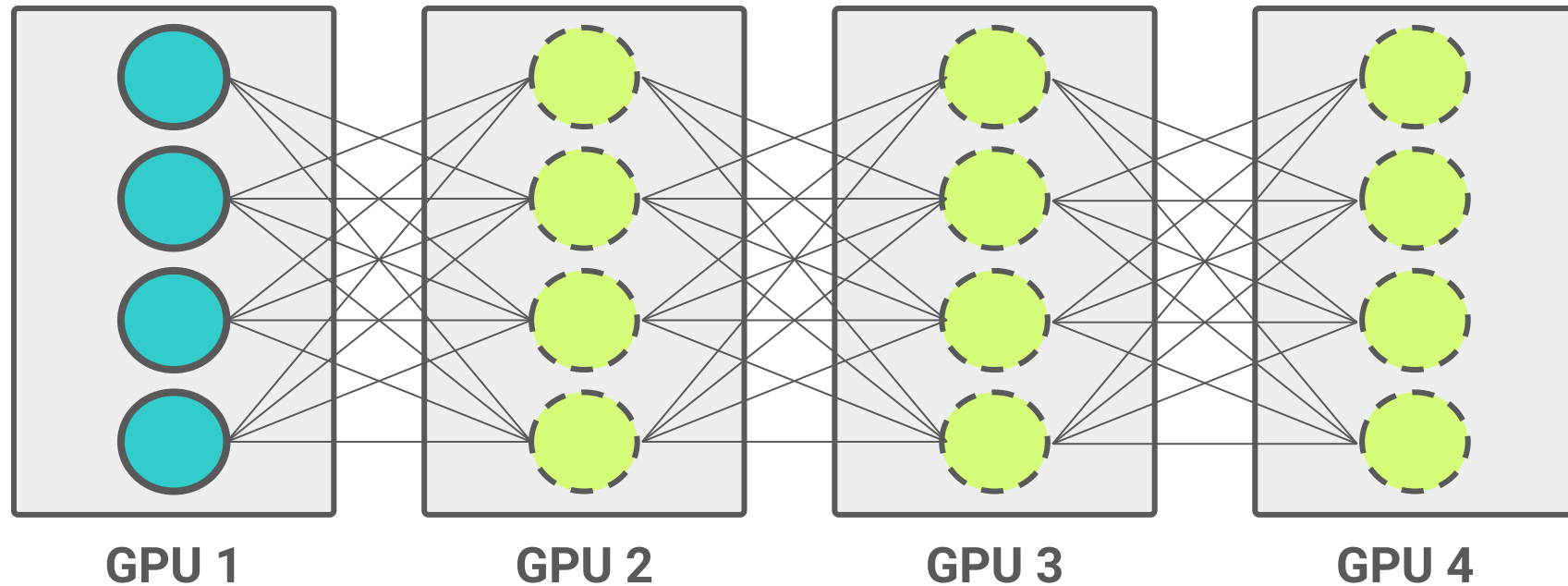GPU 1                GPU 2                GPU 3                GPU 4

Idle                 Idle

# Pipeline Parallelism



**key idea:** split mini-batch into sequential micro-batches
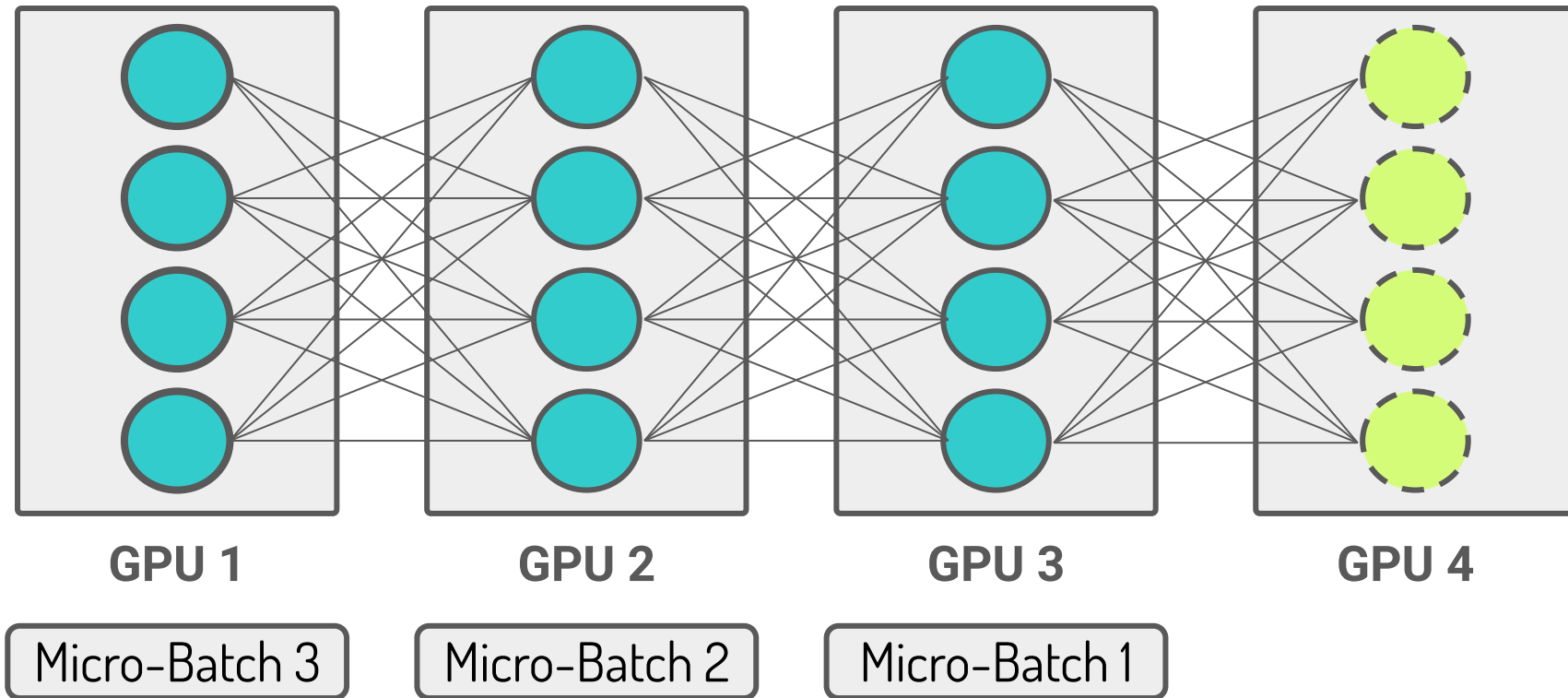
# Pipeline Parallelism



GPU 1          GPU 2          GPU 3          GPU 4

Micro-Batch 1

# Pipeline Parallelism



GPU 1          GPU 2          GPU 3          GPU 4

Micro-Batch 2     Micro-Batch 1

# Pipeline Parallelism



GPU 1　　　　GPU 2　　　　GPU 3　　　　GPU 4

Micro-Batch 3　　Micro-Batch 2　　Micro-Batch 1

14

# Pipeline Parallelism



| GPU 1 | GPU 2 | GPU 3 | GPU 4 |

| Micro-Batch 4 | Micro-Batch 3 | Micro-Batch 2 | Micro-Batch 1 |

# Pipeline Parallelism



GPU 1

GPU 2

GPU 3

GPU 4

Idle

Micro-Batch 4

Micro-Batch 3

Micro-Batch 2

# Pros and cons of model parallelism

Pros:
- Can train bigger models.
- Implemented on Pytorch.

Cons:
- Not found in the distribution strategy of default libraries such as Tensorflow. (Mesh Tensorflow)
- Tricky to design an implement.

# Data Parallelism (Large Batch Training)

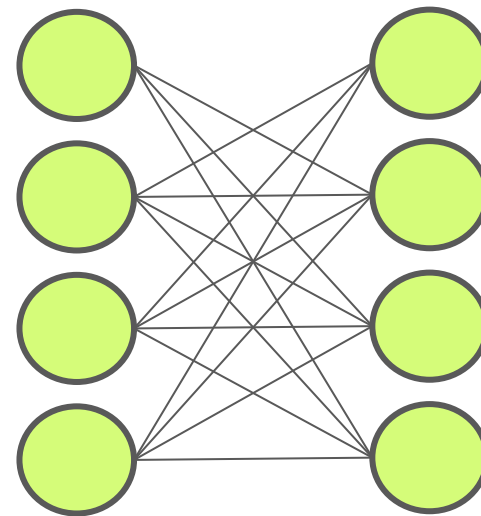**split the data across devices**

each device sees a fraction of the batch

each device replicates the model

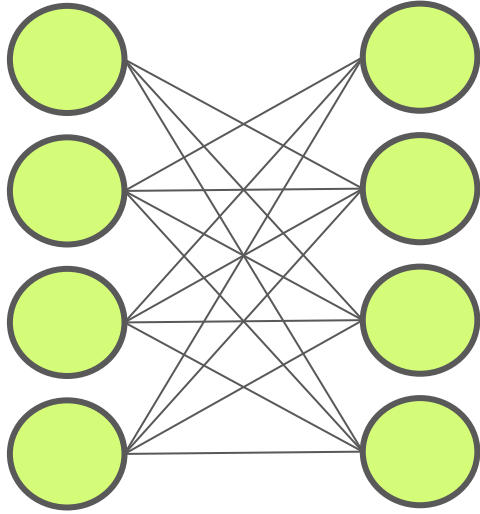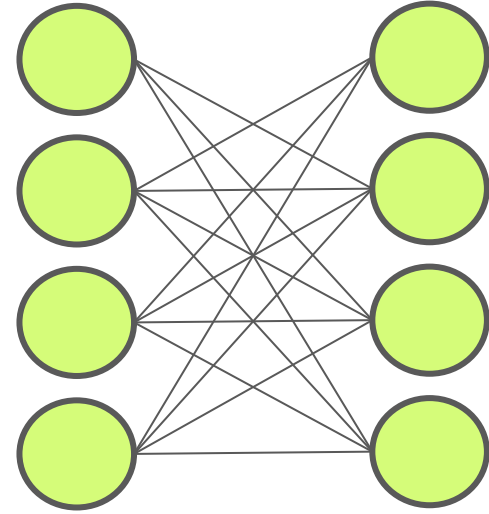each device replicates the optimizer

# Data Parallelism



GPU 1                              GPU 2

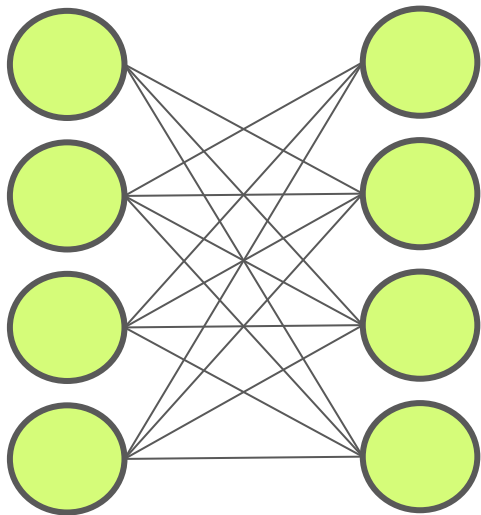GPUs could be on same or multiple nodes

# Get a batch of data



GPU 1                                                    GPU 2

Credit: Scaling Laws for Neural Language Models, stanford cs329 slides

# Split batch across devices



GPU 1

GPU 2

Credit: Scaling Laws for Neural Language Models, stanford cs329 slides

# Parallel forward passes



GPU 1

GPU 2

Credit: Scaling Laws for Neural Language Models, stanford cs329 slides

# Parallel forward passes



GPU 1

GPU 2

Credit: Scaling Laws for Neural Language Models, stanford cs329 slides

# Backpropagate gradients



GPU 1

GPU 2

Credit: Scaling Laws for Neural Language Models, stanford cs329 slides
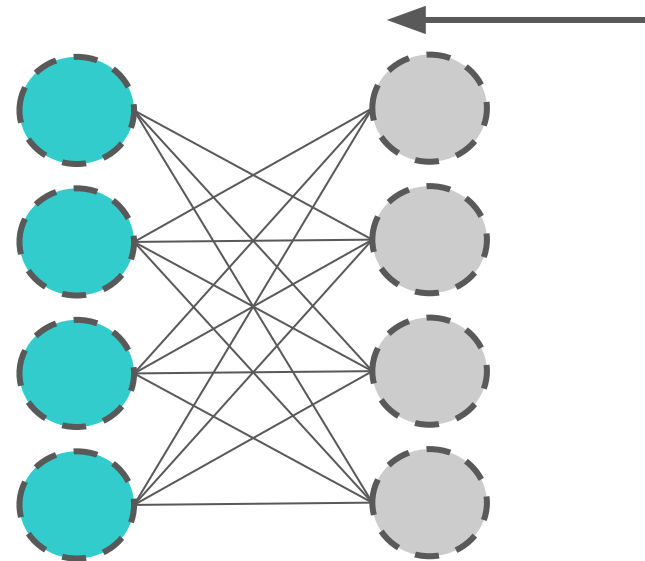
# Backpropagate gradients



GPU 1

Gradients GPU1

GPU 2

Gradients GPU2

Credit: Scaling Laws for Neural Language Models, stanford cs329 slides
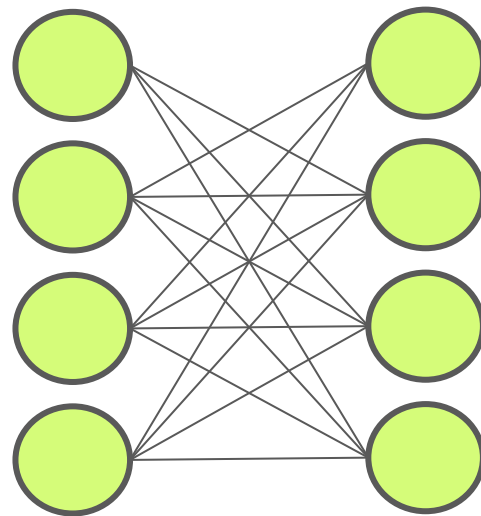
# Share gradients among GPUs and update



GPU 1                                    GPU 2

All reduce operation
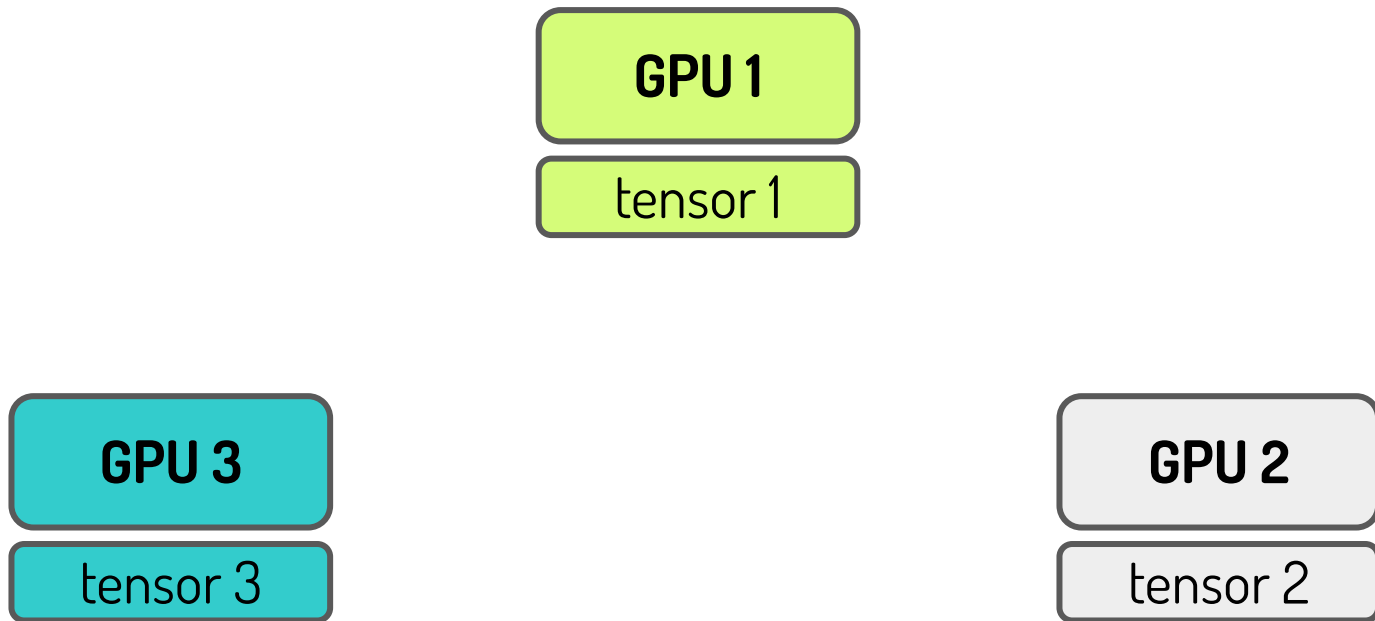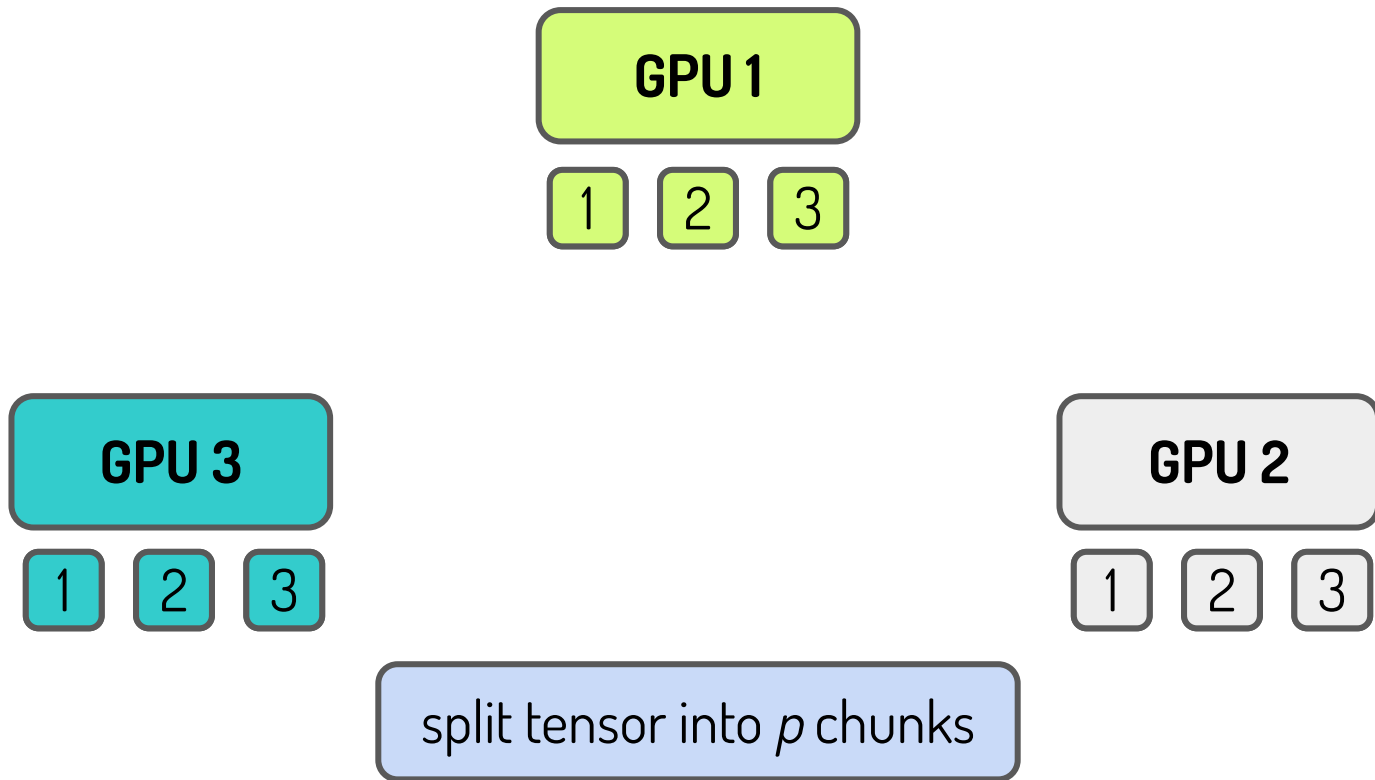
# Start the next step with a new minibatch of data

GPU 1

GPU 2

**all parameters stay synchronized!**
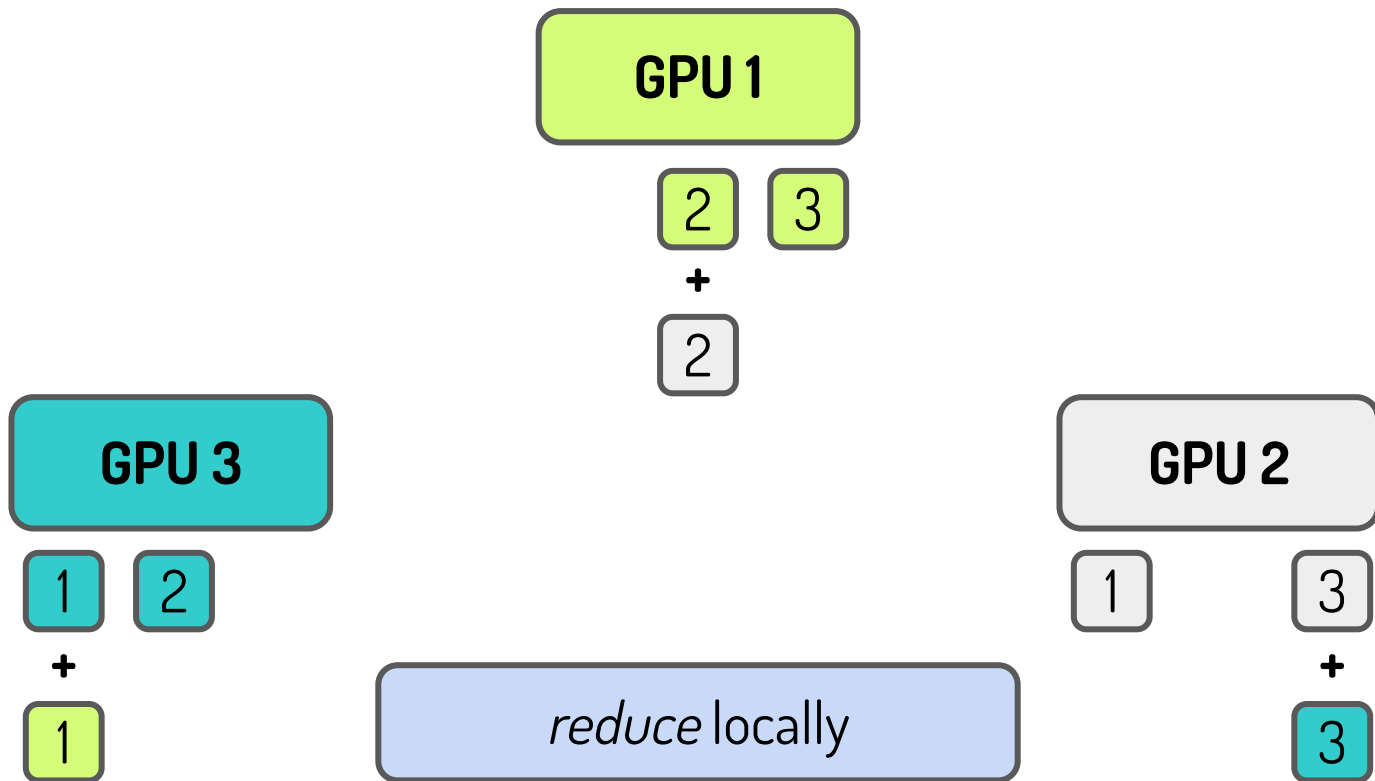
# So, what's All-Reduce?

**GPU 1**

tensor 1

**GPU 3**

tensor 3

**GPU 2**

tensor 2
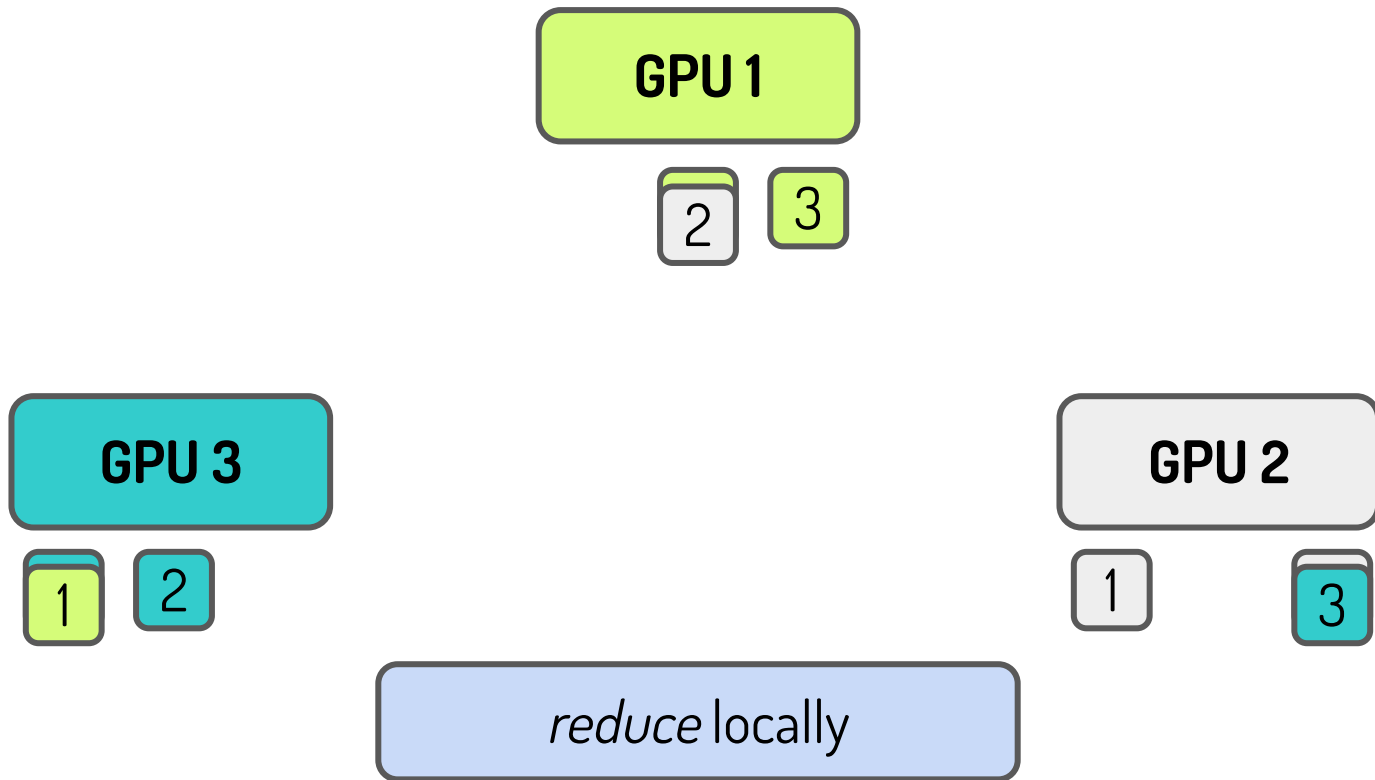
# Ring All-Reduce

# Ring All-Reduce

# Ring All-Reduce



GPU 1

2  3
+
2

GPU 3

1  2
+
1

GPU 2

1  3
+
3

*reduce* locally

# Ring All-Reduce

# Ring All-Reduce



GPU 1

2  3

GPU 3

1  2

GPU 2

1  3

*share* reduced chunk

# Ring All-Reduce



GPU 1

3
+
3

GPU 3

2
+
2

GPU 2

1
+
1

*reduce* locally

# Ring All-Reduce

**GPU 1**

3

**GPU 3**

2

**GPU 2**

1

everyone has a chunk of the result after *(p-1)* share-reduce steps

# Ring All-Reduce



GPU 1

3

GPU 3

GPU 2

2

1

*share* reduced chunk

# Ring All-Reduce
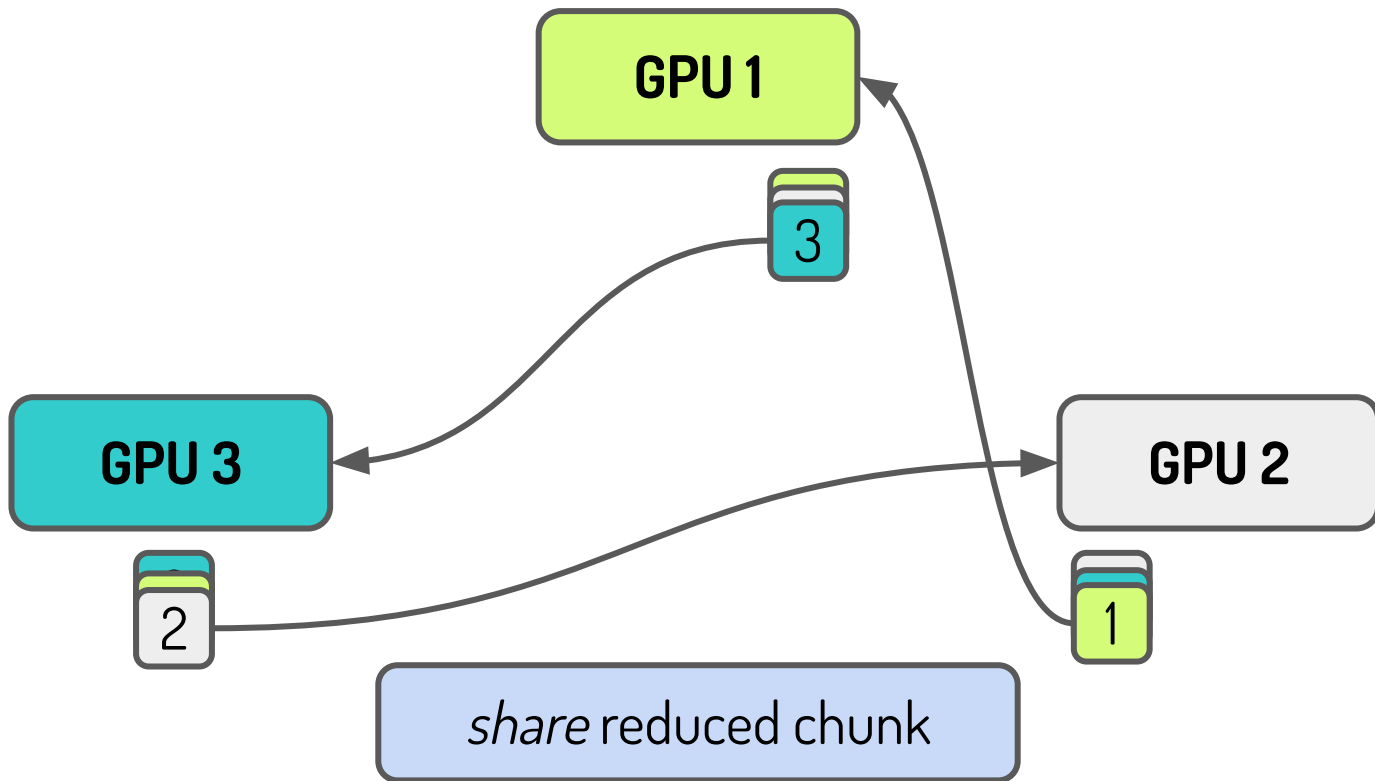
# Ring All-Reduce



GPU 1

3

GPU 3

2    3

GPU 2

1    2

*share* reduced chunk

# Ring All-Reduce

# Ring All-Reduce Advantages

Naive method = p senders x (p – 1) receivers x o(n) tensor = **o(np$^2$)**
everyone does **o(np)** work.

Manager node method = (p-1) x 2 transfers x o(n) tensor = **o(np)**
manager does **o(np)** work

Ring All-Reduce = p senders x 1 receiver x o(n/p) tensor x (p-1) rounds x 2
phases  = **o(np)**
everyone does <u>equal</u> **o(n)** work (independent of p)

# Before the demo, what strategies are implemented in Tensorflow

| Training API | MirroredStrategy | TPUStrategy | MultiWorkerMirroredStrategy | CentralStorageStrategy | ParameterServerStrategy |
|---|---|---|---|---|---|
| Keras API | Supported | Supported | Supported | Experimental support | Supported planned post 2.4 |
| Custom training loop | Supported | Supported | Supported | Experimental support | Experimental support |
| Estimator API | Limited Support | Not supported | Limited Support | Limited Support | Limited Support |

# Demos

Using data parallelism with model.fit

Using data parallelism with a custom training loop.

(https://github.com/juanma9613/Reducing-deep-learning-training-times-Pycon2021)

# THANK YOU!

- https://www.linkedin.com/in/juan-manuel-munoz-betancur/
- jmunozb@eafit.edu.co

Link to demo:

- https://github.com/juanma9613/Reducing-deep-learning-training-times-Pycon2021

**Factored is hiring!**