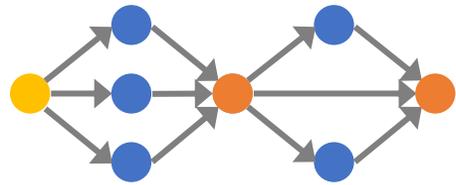# 15-442/15-642: Machine Learning Systems

# Parallelization Part 2
# (Model and Pipeline Parallelism)

**Tianqi Chen and Zhihao Jia**

Carnegie Mellon University
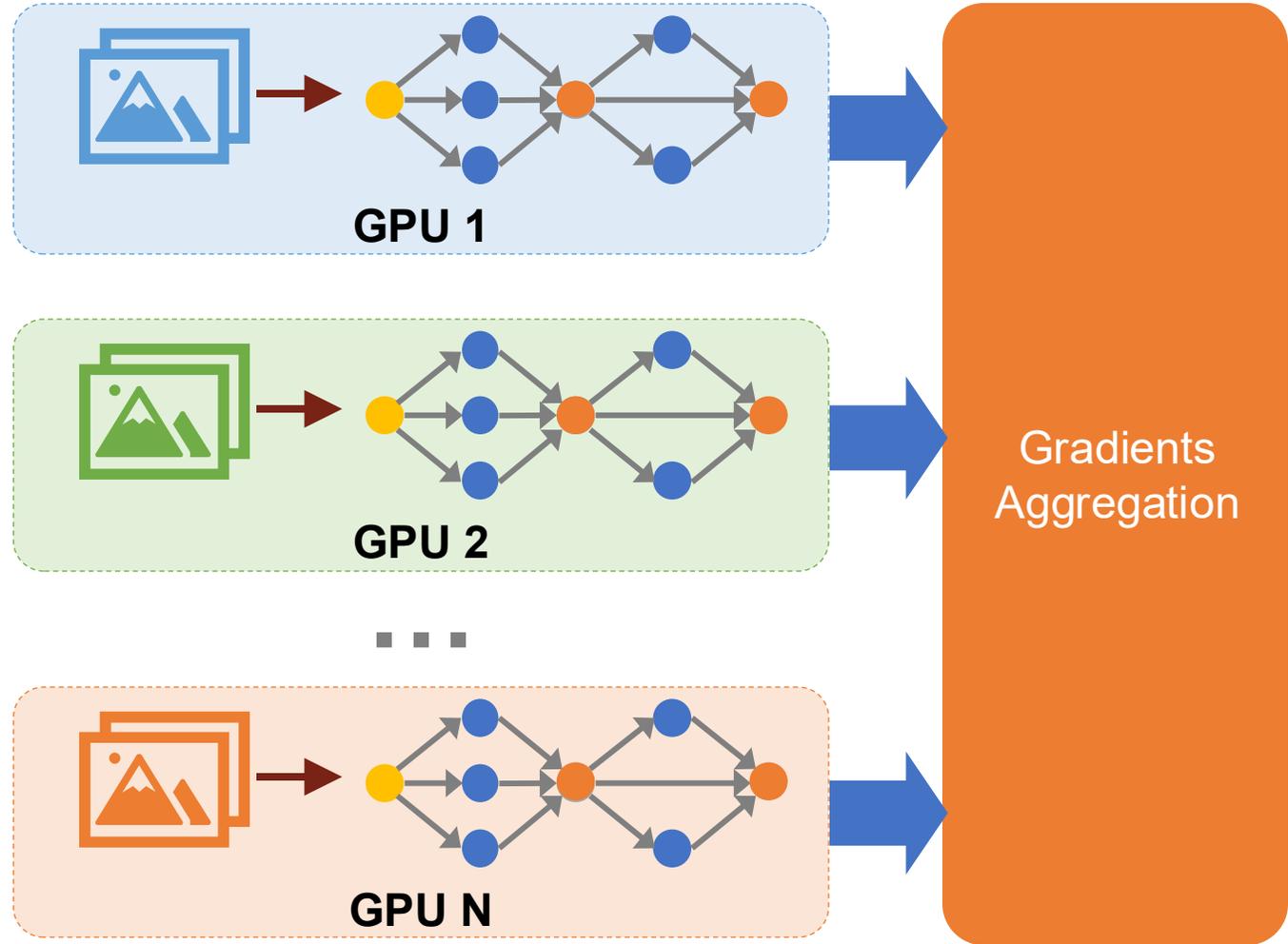
# Recap: Data Parallelism



**ML Model**

**Training Dataset**

$$w_i := w_i - \gamma \nabla L(w_i) = w_i - \frac{\gamma}{n} \sum_{j=1}^{n} \nabla L_j(w_i)$$

**GPU 1**

**GPU 2**

**GPU N**

Gradients Aggregation

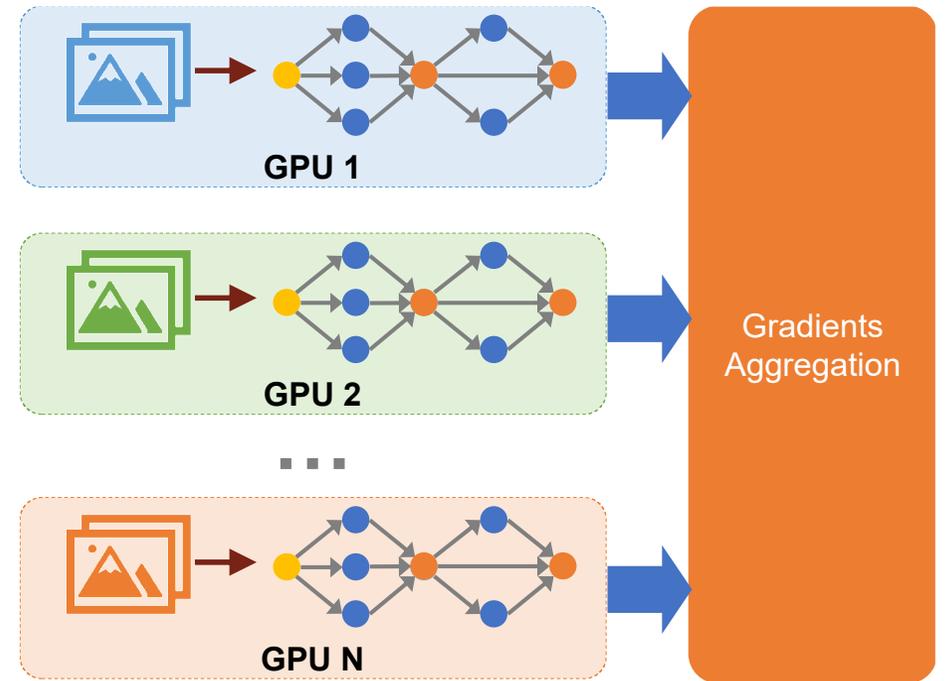1. Partition training data into batches

2. Compute the gradients of each batch on a GPU
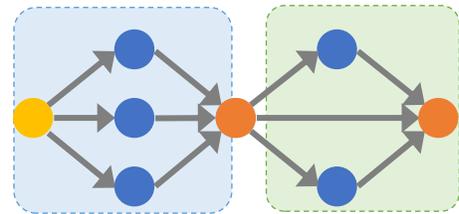
3. Aggregate gradients across GPUs

# Recap: An Issue with Data Parallelism

- Each GPU saves a replica of the entire model

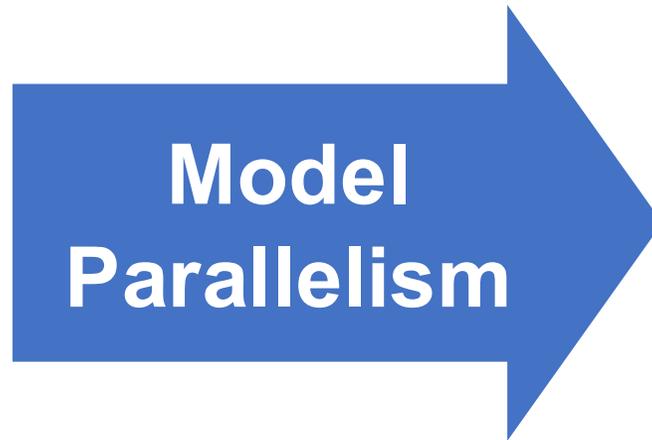- Cannot train large models that exceed GPU device memory

# Model Parallelism

- Split a model into multiple subgraphs and assign them to different devices
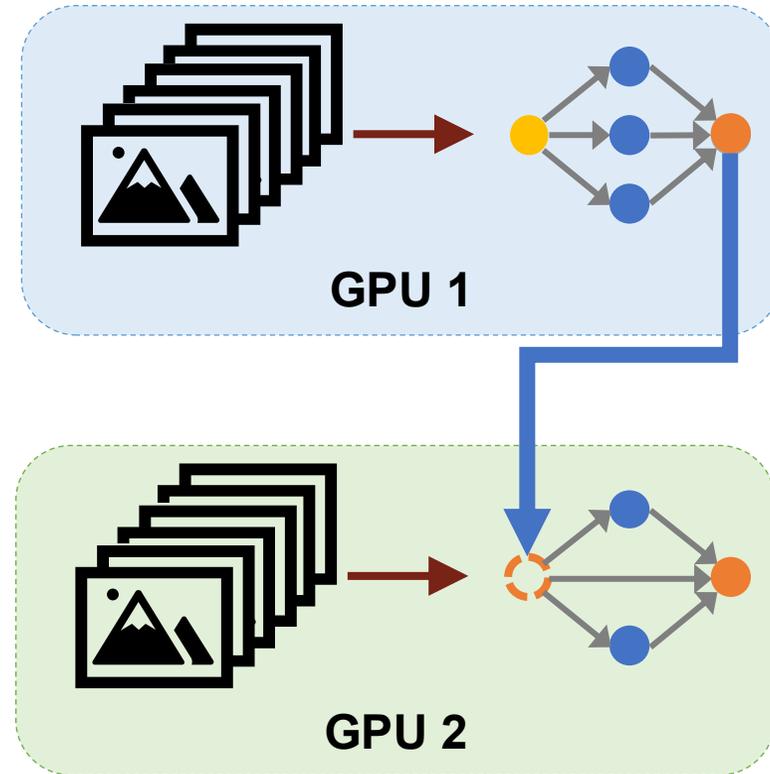
**ML Model**

**Training Dataset**

**Model Parallelism**

**GPU 1**

**GPU 2**

Transfer intermediate results between devices

$$w_i := w_i - \gamma \nabla L(w_i) = w_i - \frac{\gamma}{n} \sum_{j=1}^{n} \nabla L_j(w_i)$$

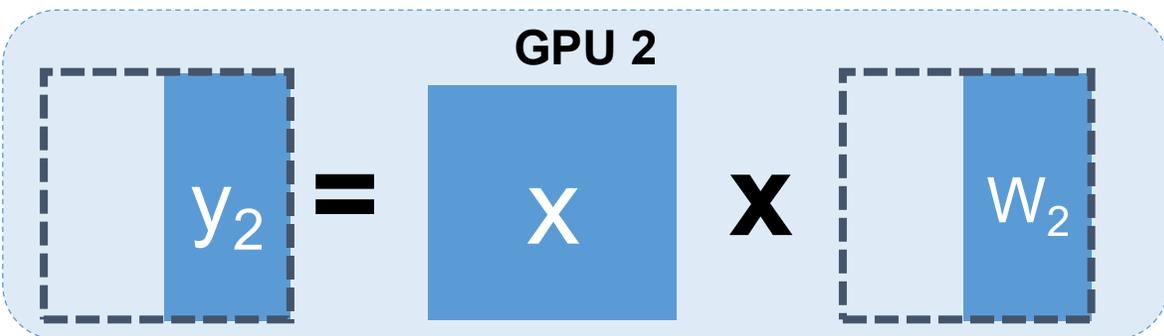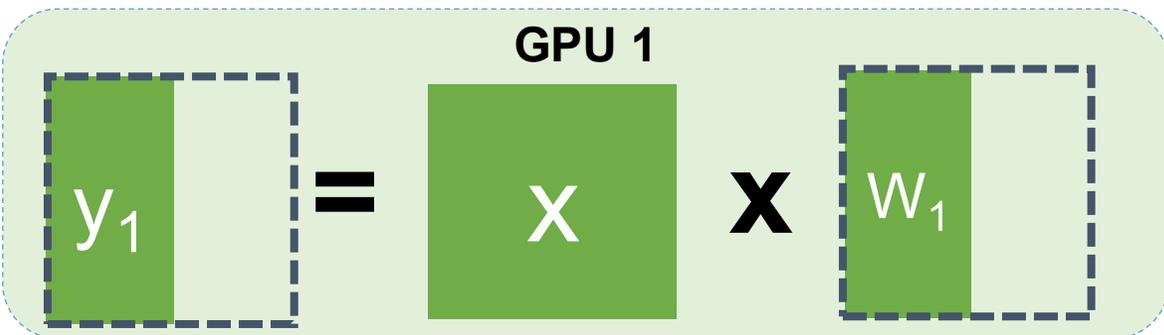# How to parallelize DNN Training?

- Data parallelism

- Model parallelism
  - Tensor model parallelism
  - Pipeline model parallelism

# Tensor Model Parallelism

$$y = x \times W$$

output = input × parameters
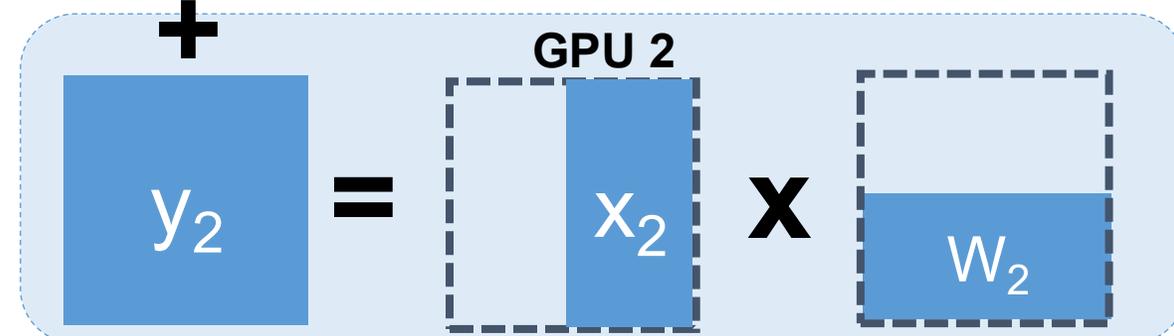
- Partition parameters/gradients *within* a layer



Tensor Model Parallelism (partition output)

Tensor Model Parallelism (reduce output)
$$y = y1 + y2$$

8

# Comparing Data and Tensor Model Parallelism



$$y = Wx$$

Data parallelism

| Forward Processing | Backward Propagation | Gradients Sync |
|---|---|---|
| 0 | 0 | $O(C_{out} * C_{in})$ |

Communication Cost of Data Parallelism

# Comparing Data and Tensor Model Parallelism



Tensor Model Parallelism (partition output)

| Forward Processing | Backward Propagation | Gradients Sync |
|:---:|:---:|:---:|
| 0 | $O(B * C_{in})$ | 0 |

Communication Cost of Tensor Model Parallelism

# Comparing Data and Tensor Model Parallelism



$$y = x \times W$$

with dimensions $C_{out}$, $B$ for $y$; $C_{in}$, $B$ for $x$; $C_{out}$, $C_{in}$ for $W$.

**GPU 1**

$$y_1 = x_1 \times W_1$$

**+**

**GPU 2**

$$y_2 = x_2 \times W_2$$

Tensor Model Parallelism (Reduce output)
$$y = y1 + y2$$

| Forward Processing | Backward Propagation | Gradients Sync |
|---|---|---|
| $O(B * C_{out})$ | 0 | 0 |

Communication Cost of Tensor Model Parallelism

11

# Comparing Data and Tensor Model Parallelism

- Data parallelism: $O(C_{out} * C_{in})$
- Tensor model parallelism (partition output): $O(B * C_{in})$
- Tensor model parallelism (reduce output): $O(B * C_{out})$

- **The best strategy depends on the model and underlying machine**

# Combine Data and Model Parallelism

# Example: Convolutional Neural Networks



Classification

Retrieval

Detection

Segmentation

Self-Driving

Synthesis

# Convolution

- Convolve the filter with the image: slide over the image spatially and compute dot products



Source pixel

$$(-1 \times 3) + (0 \times 0) + (1 \times 1) +$$
$$(-2 \times 2) + (0 \times 6) + (2 \times 2) +$$
$$(-1 \times 2) + (0 \times 4) + (1 \times 1) = -3$$

Convolution filter
(Sobel Gx)

Destination pixel

# CNNs

- A sequence of convolutional layers, interspersed by pooling, normalization, and activation functions



Low-level features → Mid-level features → High-level features → Linearly separable classifier

VGG-16 Conv1_1          VGG-16 Conv3_2          VGG-16 Conv5_3

[Zeiler and Fergus 2013]

16

# Parallelizing Convolutional Neural Networks

- Convolutional layers
  - 90-95% of the computation
  - 5% of the parameters
  - Very large intermediate activations

- Fully-connected layers
  - 5-10% of the computation
  - 95% of the parameters
  - Small intermediate activations

- **Discussion: how to parallelize CNNs?**

**Data parallelism**

**Tensor model parallelism**

# Parallelizing Convolutional Neural Networks

- Data parallelism for convolutional layers
- Tensor model parallelism for fully-connected layers

# Example: Parallelizing Transformers

- Transformer: attention mechanism for language understanding



Ashish Vaswani et. al. Attention is all you need.

# A Single Transformer Layer



20

# Parallelizing Fully-Connected Layers in Transformers

$$Y = GeLU(X \times A)$$
$$Z = Dropout(Y \times B)$$

identity layer

reduction layer



Tensor model parallelism
(partition output)

Tensor model parallelism
(reduce output)

# Multi-Head Self-Attention



$$Z_i = A(Q_i, K_i, V_i) = softmax\left(\frac{Q_i K_i^T}{\sqrt{d}}\right) V_i$$

$$Z = MultiHead(Q, K, V) = Concat(Z_0, \ldots, Z_7) W^o$$

# Parallelizing Self-Attention Layers in Transformers

$$Y_i = A(Q_i, K_i, V_i) = softmax\left(\frac{Q_i K_i^T}{\sqrt{d}}\right) V_i$$

$$Z = MultiHead(Q, K, V) = Concat(Y_0, \ldots, Y_h) W^o$$



**Parallelizing across attention heads**

**Tensor model parallelism (reduce output)**

# Parallelizing Transformers



Scale to 512 GPUs by combining data and model parallelism

# How to parallelize DNN Training?

- Data parallelism
- Model parallelism
  - Tensor model parallelism
  - **Pipeline model parallelism**

# An Issue with Model Parallelism

- Under-utilization of compute resources
- Low overall throughput due to resource utilization

# Pipeline Model Parallelism

- **Mini-batch**: the number of samples processed in each iteration

- Divide a mini-batch into multiple **micro-batches**

- Pipeline the forward and backward computations across micro-batches



**Model Parallelism**

Forward Pass | Backward Pass | Idle

**Pipeline Model Parallelism**

All inputs use weights from last flush

Pipeline flush: add gradients

Forward Pass | Backward Pass | Idle

32

# Pipeline Model Parallelism: Device Utilization

- $m$ : micro-batches in a mini-batch
- $p$: number of pipeline stages
- All stages take $t_f / t_b$ to process a forward (backward) micro-batch



$$BubbleFraction = \frac{(p-1)*(t_f+t_b)}{m*t_f + m*t_b} = \frac{p-1}{m}$$

# Improving Pipeline Parallelism Efficiency

- $m$ : number of micro-batches in a mini-batch
  - Increase mini-batch size or reduce micro-batch size
  - Caveat: large mini-batch sizes can lead to accuracy loss; small micro-batch sizes reduce GPU utilization

- $p$: number of pipeline stages
  - Decrease pipeline depth
  - Caveat: increase stage size



$$BubbleFraction = \frac{(p-1)*(t_f+t_b)}{m*t_f + m*t_b} = \frac{p-1}{m}$$

# Pipeline Model Parallelism: Memory Requirement

- An issue: we need to keep the intermediate activations of **all micro-batches** before back propagation



**Can we improve the pipeline schedule to reduce memory requirement?**

# Pipeline Parallelism with 1F1B Schedule

- One-Forward-One-Backward in the steady state
- Limit the number of in-flight micro-batches to the pipeline depth
- **Reduce memory footprint of pipeline parallelism**
- **Doesn't reduce pipeline bubble**

**Can we reduce pipeline bubble?**



# in-flight mciro-batches = **8**

**Pipeline parallelism with GPipe's schedule**

# in-flight mciro-batches = **4**

**Pipeline parallelism with 1F1B schedule**

# Pipeline Parallelism with Interleaved 1F1B Schedule

- Further divide each stage into $v$ sub-stages

- The forward (backward) time of each sub-stage is $\frac{t_f}{v}$ $(\frac{t_b}{v})$



Each device is assigned two chunks. Dark colors show the first chunk and light colors show the second chunk.

$$BubbleFraction = \frac{(p-1) * \frac{(t_f + t_b)}{v}}{m * t_f + m * t_b} = \frac{1}{v} * \frac{p-1}{m}$$

**Reduce bubble time at the cost increased communication**

# Pipeline Parallelism with Interleaved 1F1B Schedule

**Pipeline parallelism with 1F1B Schedule**

$$BubbleFraction = \frac{p-1}{m}$$

**Pipeline parallelism with interleaved 1F1B Schedule**

$$BubbleFraction = \frac{1}{v} * \frac{p-1}{m}$$



Assign multiple stages to each device

Forward Pass    Backward Pass

38

# Summary: Comparing Data/Tensor Model/Pipeline Model Parallelism



| | Data Parallelism | Tensor Model Parallelism | Pipeline Model Parallelism |
|---|---|---|---|
| **Pros** | ✓ Massively parallelizable<br>✓ Require no communication during forward/backward | ✓ Support training large models<br>✓ Efficient for models with large numbers of parameters | ✓ Support large-batch training<br>✓ Efficient for deep models |
| **Cons** | ❖ Do not work for models that cannot fit on a GPU<br>❖ Do not scale for models with large numbers of parameters | ❖ Limited parallelizability; cannot scale to large numbers of GPUs<br>❖ Need to transfer intermediate results in forward/backward | ❖ Limited utilization: bubbles in forward/backward |

# Summary: Comparing Data/Tensor Model/Pipeline Model Parallelism



**Training large models requires combining data/model/pipeline and other parallelization techniques**

|      | Data Parallelism | Model Parallelism | Pipeline Parallelism |
|------|------------------|-------------------|----------------------|
| Pros | ✓ Massively parallelizable<br>✓ Require no communication during forward/backward | ✓ Support training large models<br>✓ Efficient for models with large numbers of parameters | ✓ Support large-batch training<br>✓ Efficient for deep models |
| Cons | ❖ Do not work for models that cannot fit on a GPU<br>❖ Do not scale for models with large numbers of parameters | ❖ Limited parallelizability; cannot scale to large numbers of GPUs<br>❖ Need to transfer intermediate results in forward/backward | ❖ Limited utilization: bubbles in forward/backward |

# Example: 3D parallelism in DeepSpeed